

رسالة محمد

Deep Learning

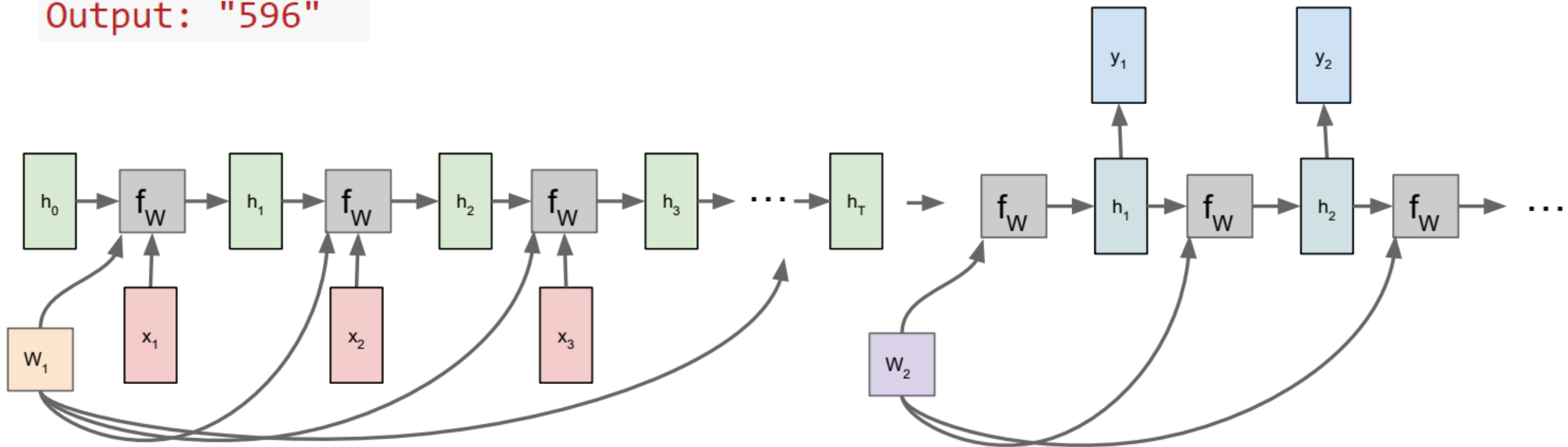
Mohammad Reza Mohammadi

2021

Sequence to Sequence

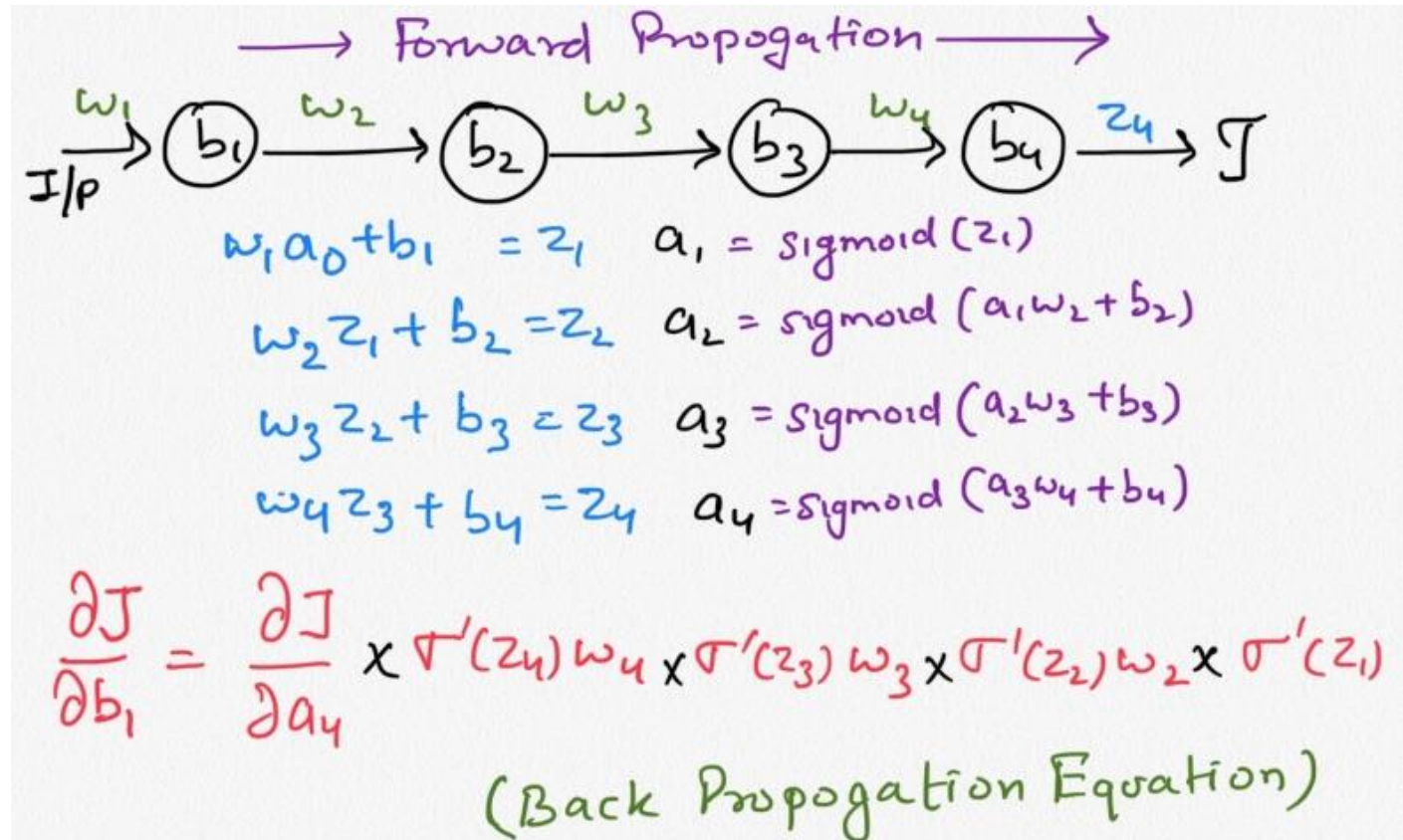
- Example: train a model to learn to add two numbers, provided as strings

Input: "535+61"
Output: "596"



Vanishing and exploding gradients

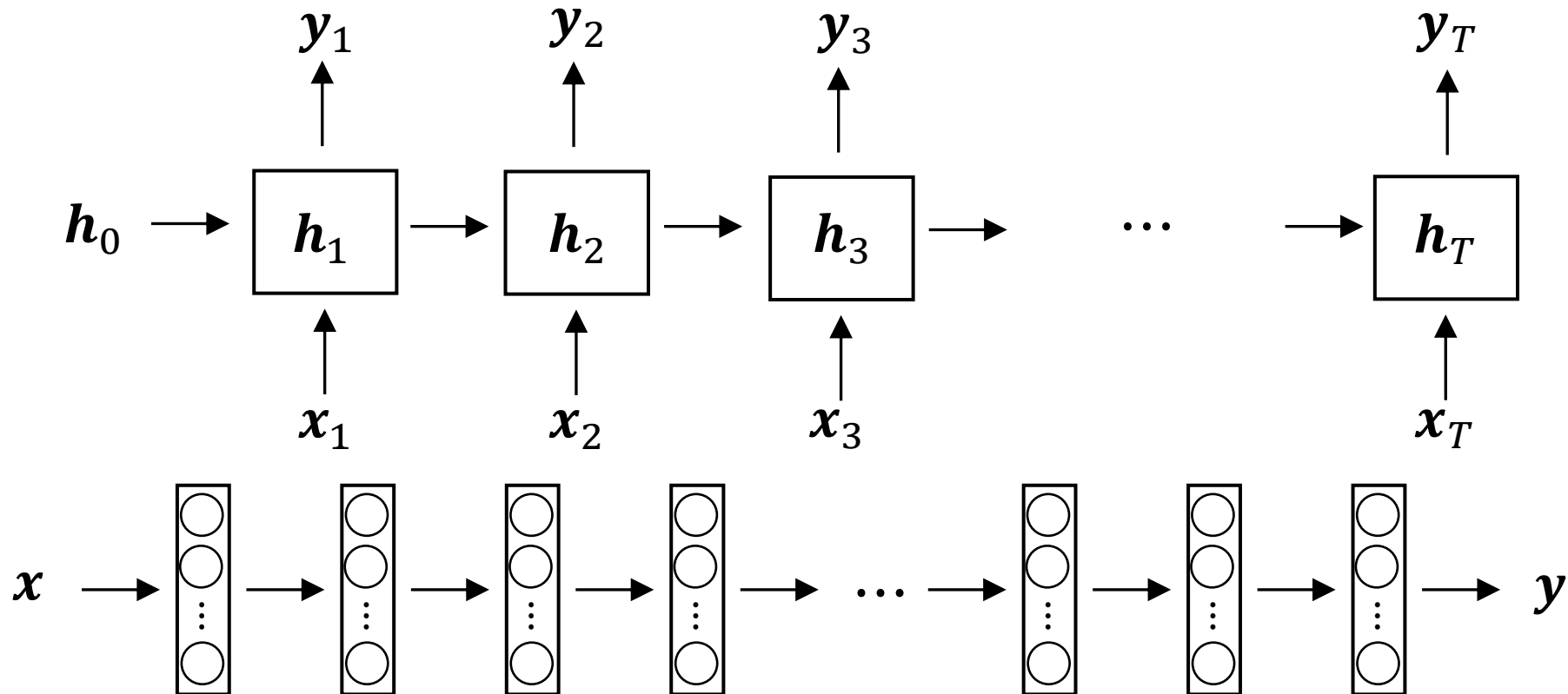
- We can avoid exploding gradients with gradient clipping



Vanishing gradients with RNNs

The **cat**, which already ate ..., **was** full.

The **cats**, which already ate ..., **were** full.



Gated RNNs

- Recurrent Neural Networks suffer from short-term memory
 - RNN's may leave out important information from the beginning
- Gated RNNs are based on the idea of creating paths through time that have derivatives that neither vanish nor explode

The **cat**, which already ate ..., **was** full.

The **cats**, which already ate ..., **were** full.

Customers Review 2,491

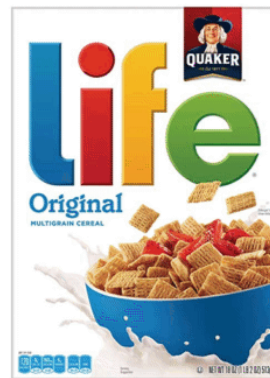


Thanos

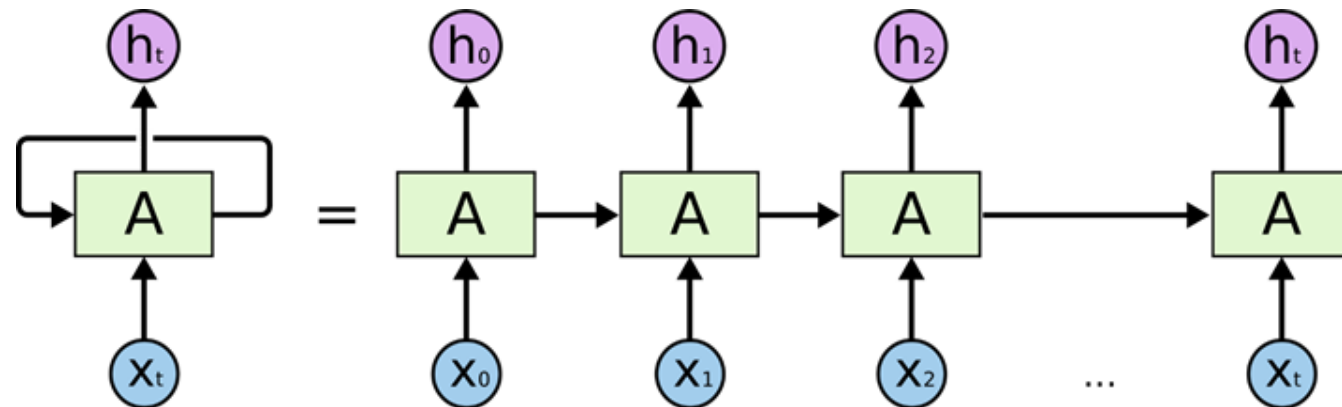
September 2018

Verified Purchase

Amazing! This box of cereal gave me a perfectly balanced breakfast, as all things should be. I only ate half of it but will definitely be buying again!



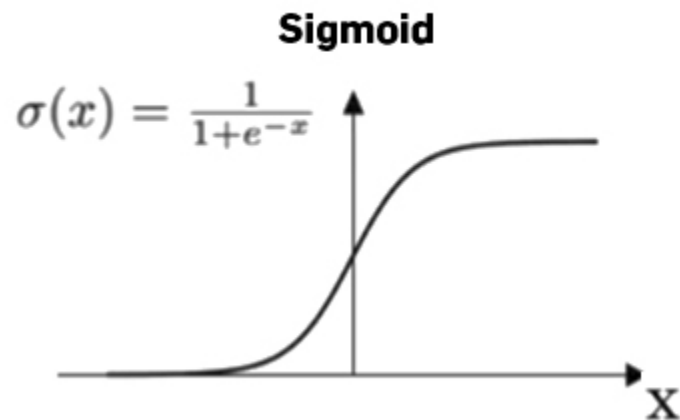
A Box of Cereal
\$3.99



Gated Recurrent Units

Simple RNN

$$\mathbf{h}^{(t)} = \tanh(\mathbf{W}_{hh} \mathbf{h}^{(t-1)} + \mathbf{W}_{xh} \mathbf{x}^{(t)} + \mathbf{b}_h)$$



GRU (simplified)

$$\tilde{\mathbf{h}}^{(t)} = \tanh(\mathbf{W}_{hh} \mathbf{h}^{(t-1)} + \mathbf{W}_{xh} \mathbf{x}^{(t)} + \mathbf{b}_h)$$

$$\mathbf{h}^{(t)} = \mathbf{u}^{(t)} \cdot \tilde{\mathbf{h}}^{(t)} + (1 - \mathbf{u}^{(t)}) \cdot \mathbf{h}^{(t-1)}$$

$$\mathbf{u}^{(t)} = \sigma(\mathbf{W}_{hu} \mathbf{h}^{(t-1)} + \mathbf{W}_{xu} \mathbf{x}^{(t)} + \mathbf{b}_u)$$

The cat, which already ate ..., was full.

Gated Recurrent Units

GRU

$$\tilde{\mathbf{h}}^{(t)} = \tanh(\mathbf{W}_{hh}(\mathbf{r}^{(t)} \cdot \mathbf{h}^{(t-1)}) + \mathbf{W}_{xh} \mathbf{x}^{(t)} + \mathbf{b}_h)$$

$$\mathbf{h}^{(t)} = \mathbf{u}^{(t)} \cdot \tilde{\mathbf{h}}^{(t)} + (1 - \mathbf{u}^{(t)}) \cdot \mathbf{h}^{(t-1)}$$

$$\mathbf{u}^{(t)} = \sigma(\mathbf{W}_{hu} \mathbf{h}^{(t-1)} + \mathbf{W}_{xu} \mathbf{x}^{(t)} + \mathbf{b}_u)$$

$$\mathbf{r}^{(t)} = \sigma(\mathbf{W}_{hr} \mathbf{h}^{(t-1)} + \mathbf{W}_{xr} \mathbf{x}^{(t)} + \mathbf{b}_r)$$

GRU (simplified)

$$\tilde{\mathbf{h}}^{(t)} = \tanh(\mathbf{W}_{hh} \mathbf{h}^{(t-1)} + \mathbf{W}_{xh} \mathbf{x}^{(t)} + \mathbf{b}_h)$$

$$\mathbf{h}^{(t)} = \mathbf{u}^{(t)} \cdot \tilde{\mathbf{h}}^{(t)} + (1 - \mathbf{u}^{(t)}) \cdot \mathbf{h}^{(t-1)}$$

$$\mathbf{u}^{(t)} = \sigma(\mathbf{W}_{hu} \mathbf{h}^{(t-1)} + \mathbf{W}_{xu} \mathbf{x}^{(t)} + \mathbf{b}_u)$$

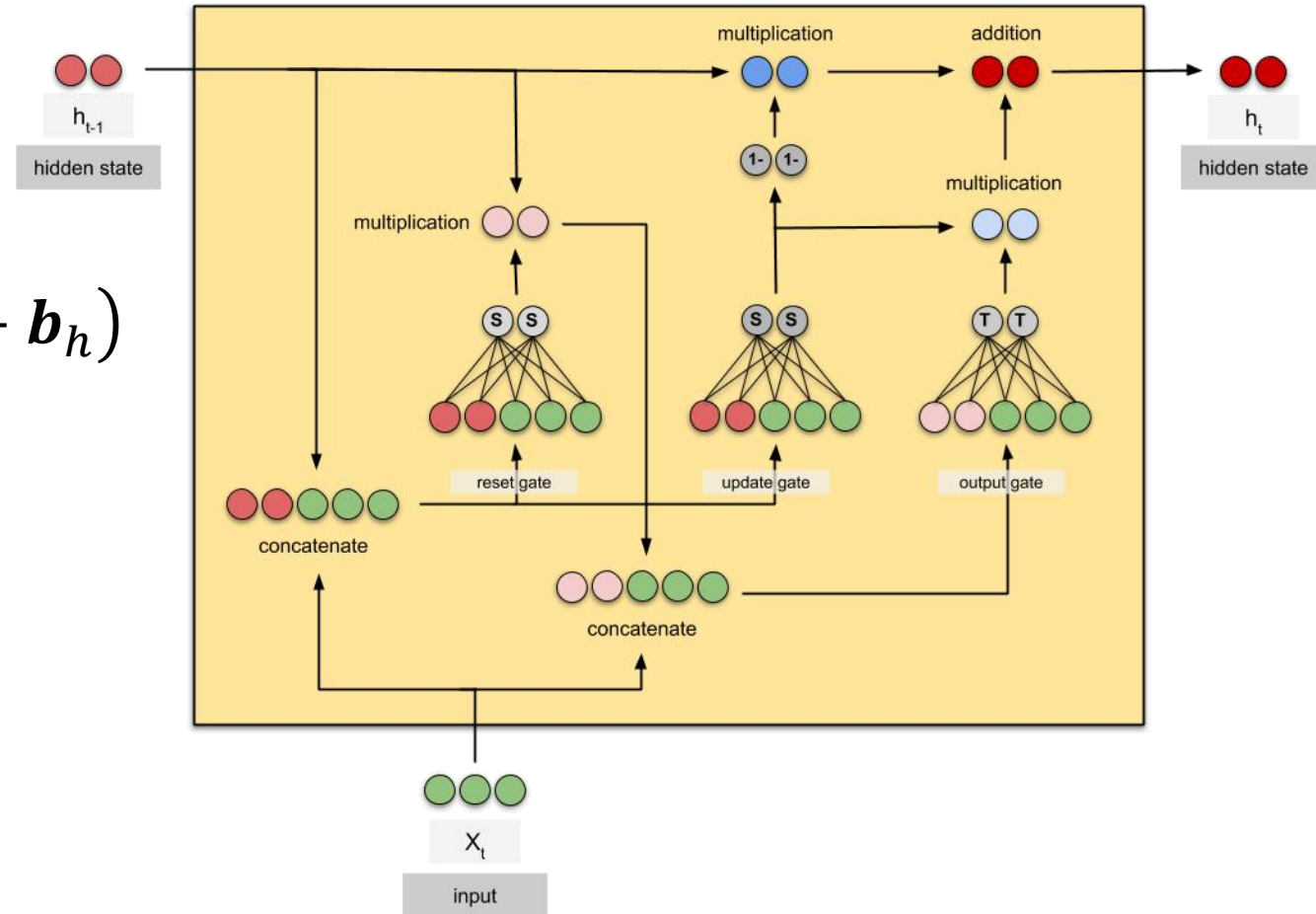
GRU

$$\tilde{\mathbf{h}}^{(t)} = \tanh(\mathbf{W}_{hh}(\mathbf{r}^{(t)} \cdot \mathbf{h}^{(t-1)}) + \mathbf{W}_{xh} \mathbf{x}^{(t)} + \mathbf{b}_h)$$

$$\mathbf{h}^{(t)} = \mathbf{u}^{(t)} \cdot \tilde{\mathbf{h}}^{(t)} + (1 - \mathbf{u}^{(t)}) \cdot \mathbf{h}^{(t-1)}$$

$$\mathbf{u}^{(t)} = \sigma(\mathbf{W}_{hu} \mathbf{h}^{(t-1)} + \mathbf{W}_{xu} \mathbf{x}^{(t)} + \mathbf{b}_u)$$

$$\mathbf{r}^{(t)} = \sigma(\mathbf{W}_{hr} \mathbf{h}^{(t-1)} + \mathbf{W}_{xr} \mathbf{x}^{(t)} + \mathbf{b}_r)$$



Long Short-Term Memory

LSTM

GRU

$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \cdot \tanh(\mathbf{c}^{(t)})$$

$$\tilde{\mathbf{h}}^{(t)} = \tanh(\mathbf{W}_{hh}(\mathbf{r}^{(t)} \cdot \mathbf{h}^{(t-1)}) + \mathbf{W}_{xh} \mathbf{x}^{(t)} + \mathbf{b}_h)$$

$$\tilde{\mathbf{c}}^{(t)} = \tanh(\mathbf{W}_{hc} \mathbf{h}^{(t-1)} + \mathbf{W}_{xc} \mathbf{x}^{(t)} + \mathbf{b}_c)$$

$$\mathbf{h}^{(t)} = \mathbf{u}^{(t)} \cdot \tilde{\mathbf{h}}^{(t)} + (1 - \mathbf{u}^{(t)}) \cdot \mathbf{h}^{(t-1)}$$

$$\mathbf{c}^{(t)} = \mathbf{u}^{(t)} \cdot \tilde{\mathbf{c}}^{(t)} + \mathbf{f}^{(t)} \cdot \mathbf{c}^{(t-1)}$$

$$\mathbf{u}^{(t)} = \sigma(\mathbf{W}_{hu} \mathbf{h}^{(t-1)} + \mathbf{W}_{xu} \mathbf{x}^{(t)} + \mathbf{b}_u)$$

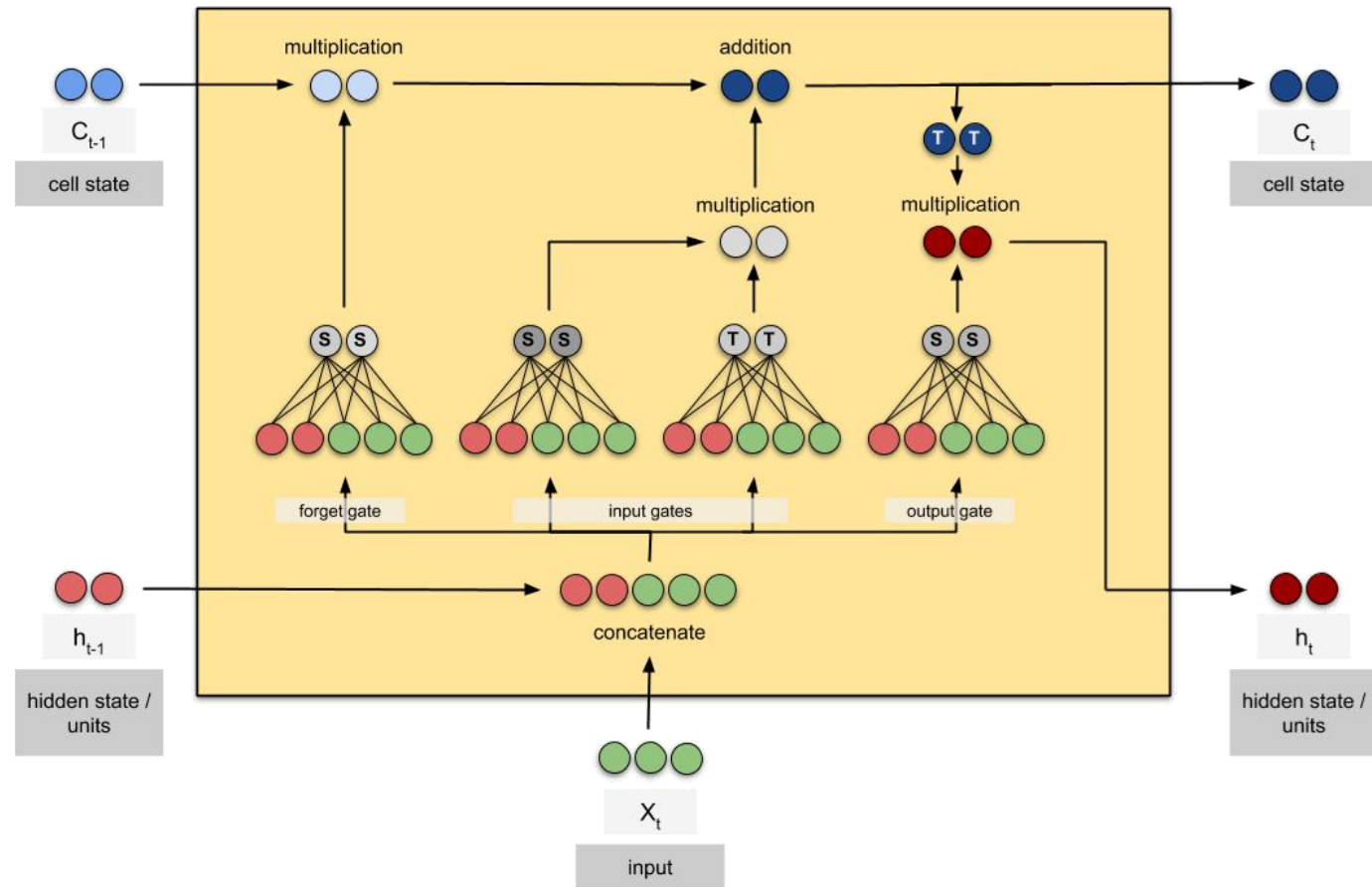
$$\mathbf{u}^{(t)} = \sigma(\mathbf{W}_{hu} \mathbf{h}^{(t-1)} + \mathbf{W}_{xu} \mathbf{x}^{(t)} + \mathbf{b}_u)$$

$$\mathbf{r}^{(t)} = \sigma(\mathbf{W}_{hr} \mathbf{h}^{(t-1)} + \mathbf{W}_{xr} \mathbf{x}^{(t)} + \mathbf{b}_r)$$

$$\mathbf{f}^{(t)} = \sigma(\mathbf{W}_{hf} \mathbf{h}^{(t-1)} + \mathbf{W}_{xf} \mathbf{x}^{(t)} + \mathbf{b}_f)$$

$$\mathbf{o}^{(t)} = \sigma(\mathbf{W}_{ho} \mathbf{h}^{(t-1)} + \mathbf{W}_{xo} \mathbf{x}^{(t)} + \mathbf{b}_o)$$

LSTM



$$h^{(t)} = o^{(t)} \cdot \tanh(c^{(t)})$$

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

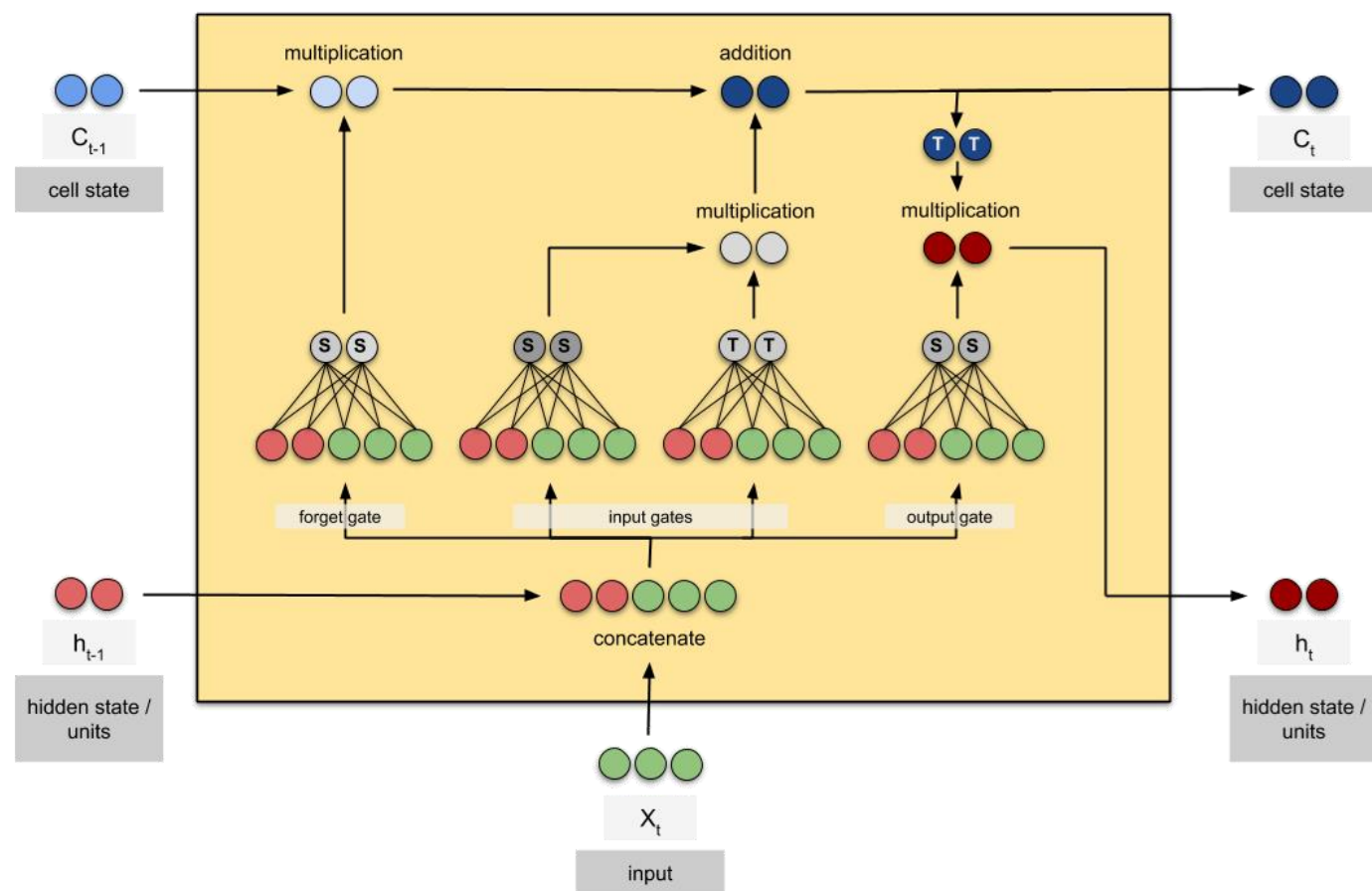
$$c^{(t)} = u^{(t)} \cdot \tilde{c}^{(t)} + f^{(t)} \cdot c^{(t-1)}$$

$$u^{(t)} = \sigma(W_{hu} h^{(t-1)} + W_{xu} x^{(t)} + b_u)$$

$$f^{(t)} = \sigma(W_{hf} h^{(t-1)} + W_{xf} x^{(t)} + b_f)$$

$$o^{(t)} = \sigma(W_{ho} h^{(t-1)} + W_{xo} x^{(t)} + b_o)$$

LSTM



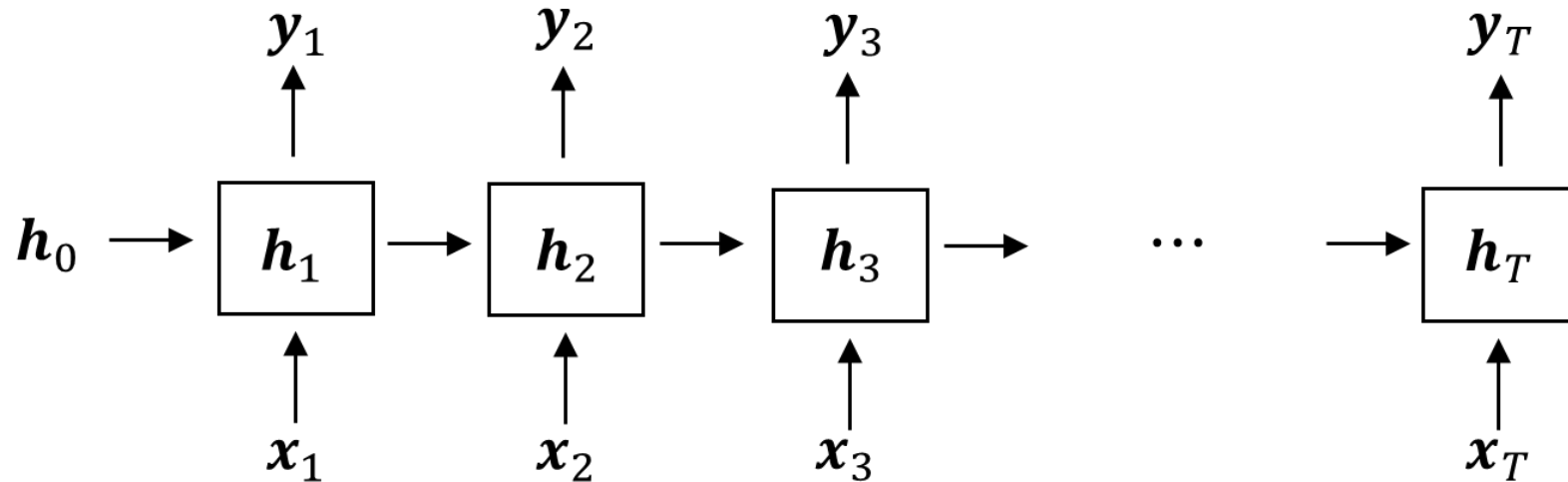
```
def LSTMCELL(prev_ct, prev_ht, input):
    combine = prev_ht + input
    ft = forget_layer(combine)
    candidate = candidate_layer(combine)
    it = input_layer(combine)
    Ct = prev_ct * ft + candidate * it
    ot = output_layer(combine)
    ht = ot * tanh(Ct)
    return ht, Ct
```

```
ct = [0, 0, 0]
ht = [0, 0, 0]
```

```
for input in inputs:
    ct, ht = LSTMCELL(ct, ht, input)
```

Bidirectional RNNs

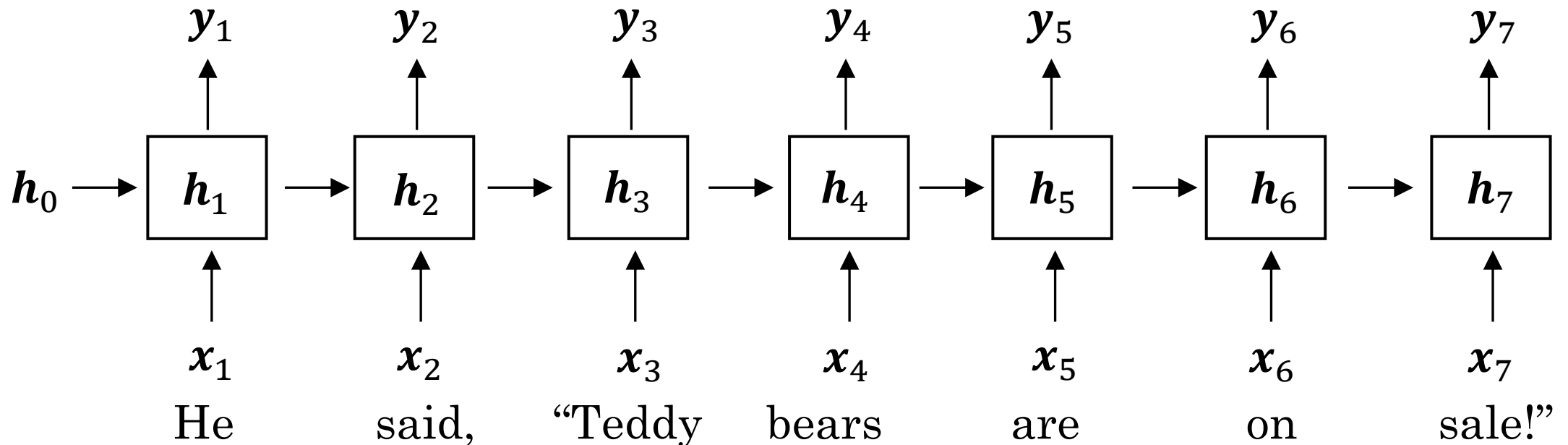
- All of the RNNs we have considered up to now have a “causal” structure
 - The state at time t only captures information from the past, $x(1), \dots, x(t-1)$, and the present input $x(t)$
- In some applications, we want to output a prediction of $y(t)$ which may depend on the whole input sequence



Bidirectional RNNs

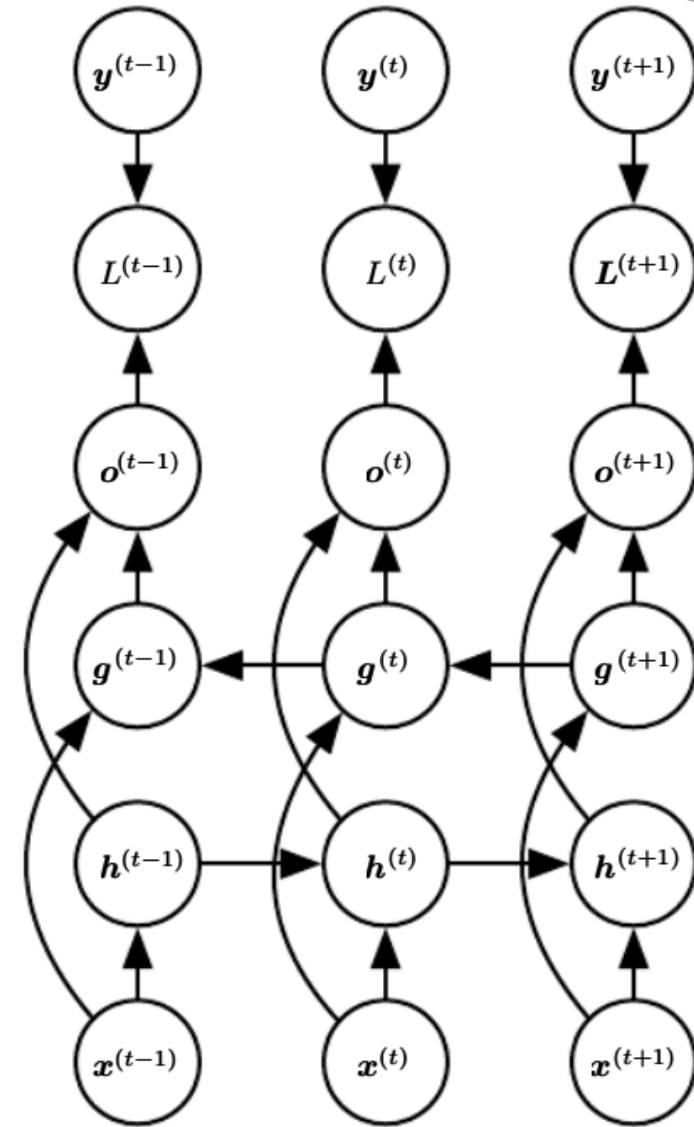
He said, “Teddy bears are on sale!”

He said, “Teddy Roosevelt was a great President!”



Bidirectional RNNs

- Bidirectional RNNs combine an RNN that moves forward through time beginning from the start of the sequence with another RNN that moves backward through time beginning from the end of the sequence
 - $h(t)$ stands for the state of the sub-RNN that moves forward through time
 - $g(t)$ stands for the state of the sub-RNN that moves backward through time
 - $o(t)$ depends on both the past and the future but is most sensitive to the input values around time t

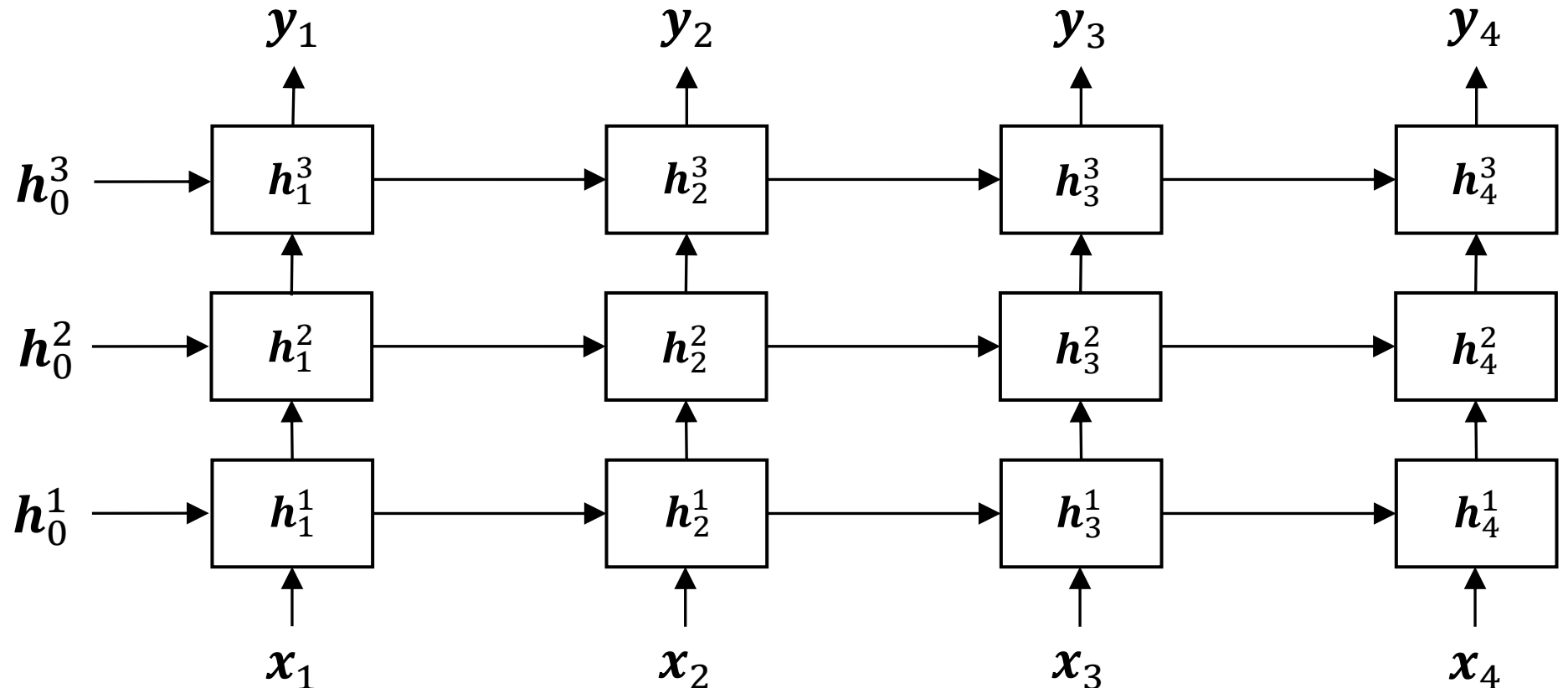
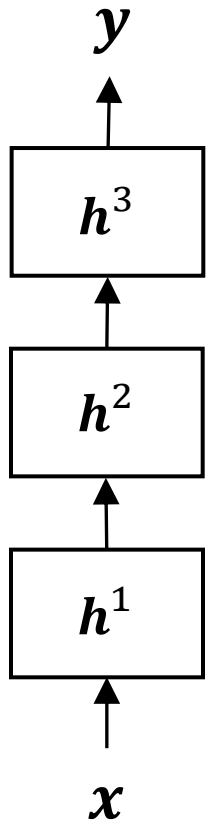


Deep Recurrent Networks

- For learning very complex functions sometimes is useful to stack multiple layers of RNNs together to build even deeper versions of these models

Deep Recurrent Networks

$$h_3^2 = \tanh(W^2 [h_2^2 \quad h_3^1] + b^2)$$



Deep Recurrent Networks

- The same parameters W^2 and b^2 are used for every one of the computations at this layer
- The first layer would have its own parameters and so on
- For RNNs, having three layers is already quite a lot.
- The blocks can be SimpleRNN, GRU, or LSTM
- They can be bidirectional

