

# Recurrent Neural Networks (RNN)

David I. Inouye



Elmore Family School of Electrical  
and Computer Engineering

Sequential data  
is natural in  
many  
applications



Text analysis



Speech recognition



Medical time series



Stock prices

# Windowing is a simple approach to handle sequential data with standard NNs

- Break up sequence into multiple fixed-length sequences:

$$\text{split}(x, y) = \{(w_j, y_j)\}_{j=1}^{\dim(x)/W}$$

- “This is a great movie.”  $\rightarrow$  { (“This is”, +1), (“is a”, +1), (“a great”, +1), (“great movie”, +1) }
- “Hello world!”  $\rightarrow$  { (“###”, “H”), (“##H”, “e”), (“#He”, “l”), (“Hel”, “l”), (“ell”, “o”), ... }
- Apply model  $f$  to each window

$$\forall (w_j, y_j) \in \text{split}(x, y), \quad \hat{y}_j = f(w_j)$$

- *Training*: Compute loss on each term or an aggregate term

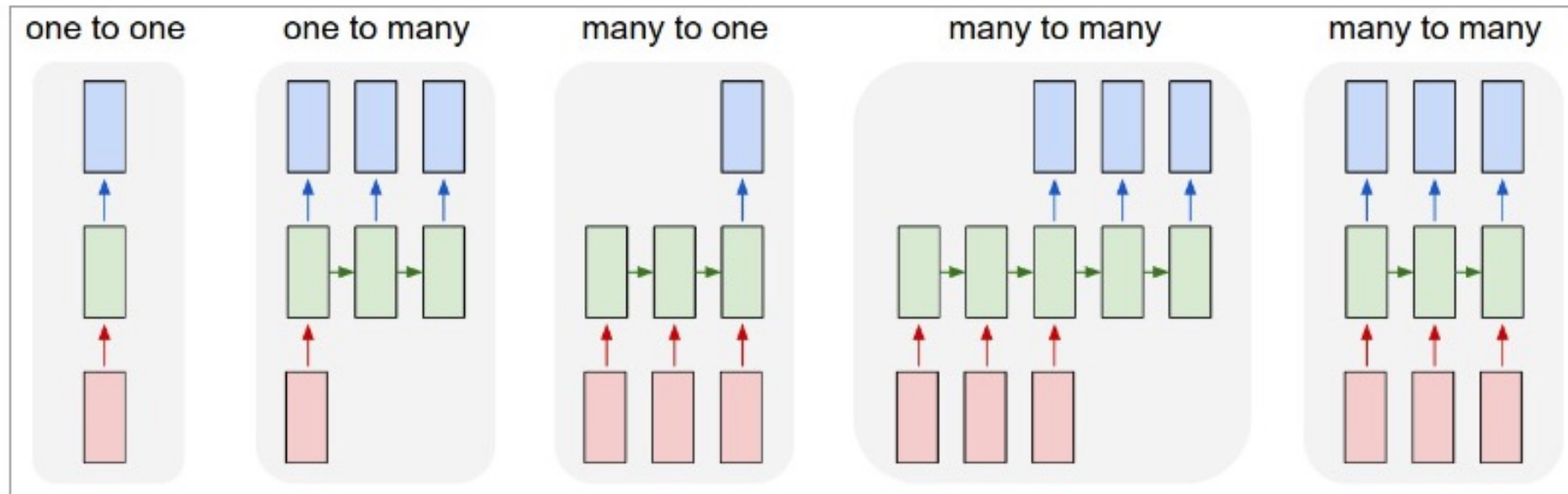
$$\sum_j \ell(\hat{y}_j, y_j) \quad \text{or} \quad \ell\left(\sum_j \hat{y}_j, y\right)$$

- *Test-time*: Concatenate or average predictions for all windows.

# While a good baseline for sequences, the windowing approach has several issues

- Fixed-window size
  - How do you choose the fixed-window size?
  - If too big, computational cost is high and learning could be slow.
  - If too small, the window may lack sufficient history to predict.
- Lacks long-range dependencies (limited to window)
  - Cannot model dependencies beyond the window size
- Predictions on each window are assumed to be independent
  - Window overlap can help as the inputs are implicitly dependent
  - Yet the outputs are not explicitly dependent

# Recurrent neural networks (RNNs) process data **sequentially** and can handle **variable-sized** input/output sequences



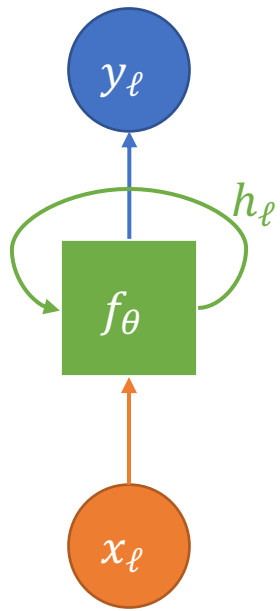
Each rectangle is a vector and arrows represent functions (e.g. matrix multiply). Input vectors are in red, output vectors are in blue and green vectors hold the RNN's state (more on this soon). From left to right: **(1)** Vanilla mode of processing without RNN, from fixed-sized input to fixed-sized output (e.g. image classification). **(2)** Sequence output (e.g. image captioning takes an image and outputs a sentence of words). **(3)** Sequence input (e.g. sentiment analysis where a given sentence is classified as expressing positive or negative sentiment). **(4)** Sequence input and sequence output (e.g. Machine Translation: an RNN reads a sentence in English and then outputs a sentence in French). **(5)** Synced sequence input and output (e.g. video classification where we wish to label each frame of the video). Notice that in every case are no pre-specified constraints on the lengths sequences because the recurrent transformation (green) is fixed and can be applied as many times as we like.

# RNNs take an input + old hidden state and produce output + new hidden state

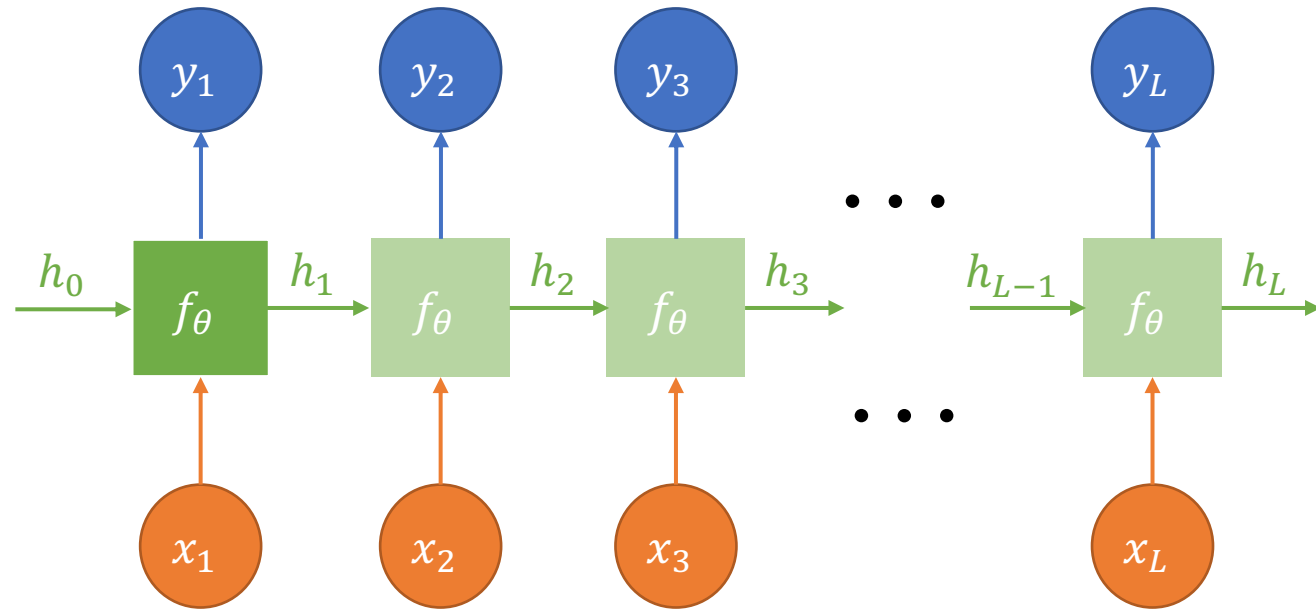
- Let  $x$ ,  $y$ , and  $h$  denote the input, output, and hidden state sequences
  - Each element in sequence could be any format including a vector, a discrete integer, or even a full tensor itself (e.g., video processing)
- Let  $L$  corresponds to length of the sequence
  - For example,  $x = (x_1, x_2, \dots, x_\ell, \dots, x_L)$
  - Note this can be different for each sample
  - For one-to-many or many-to-one, the sequences can be padded to be the same length
- RNN (parametrized by  $\theta$ ) written as recursion, where  $z_0$  is initialized to some default value:

$$y_\ell, h_\ell = f_\theta(x_\ell, h_{\ell-1})$$

# RNNs can be visualized with loop arrows or unrolled with model copies



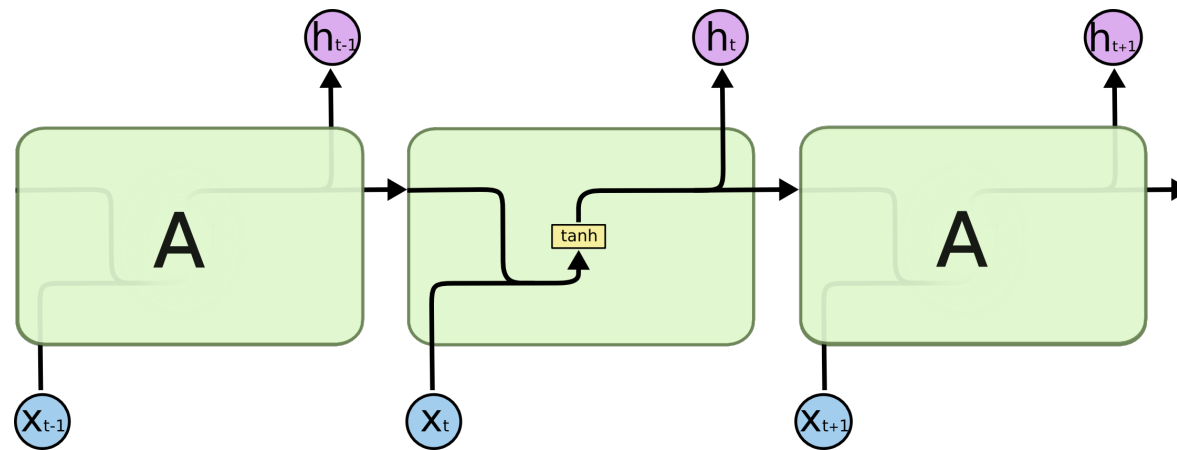
Compact recursive visualization shows the RNNs simple form.



Unrolled visualization shows that the same network is used multiple times but with different inputs and different hidden states.

# A vanilla RNN can be made with linear and activation layers

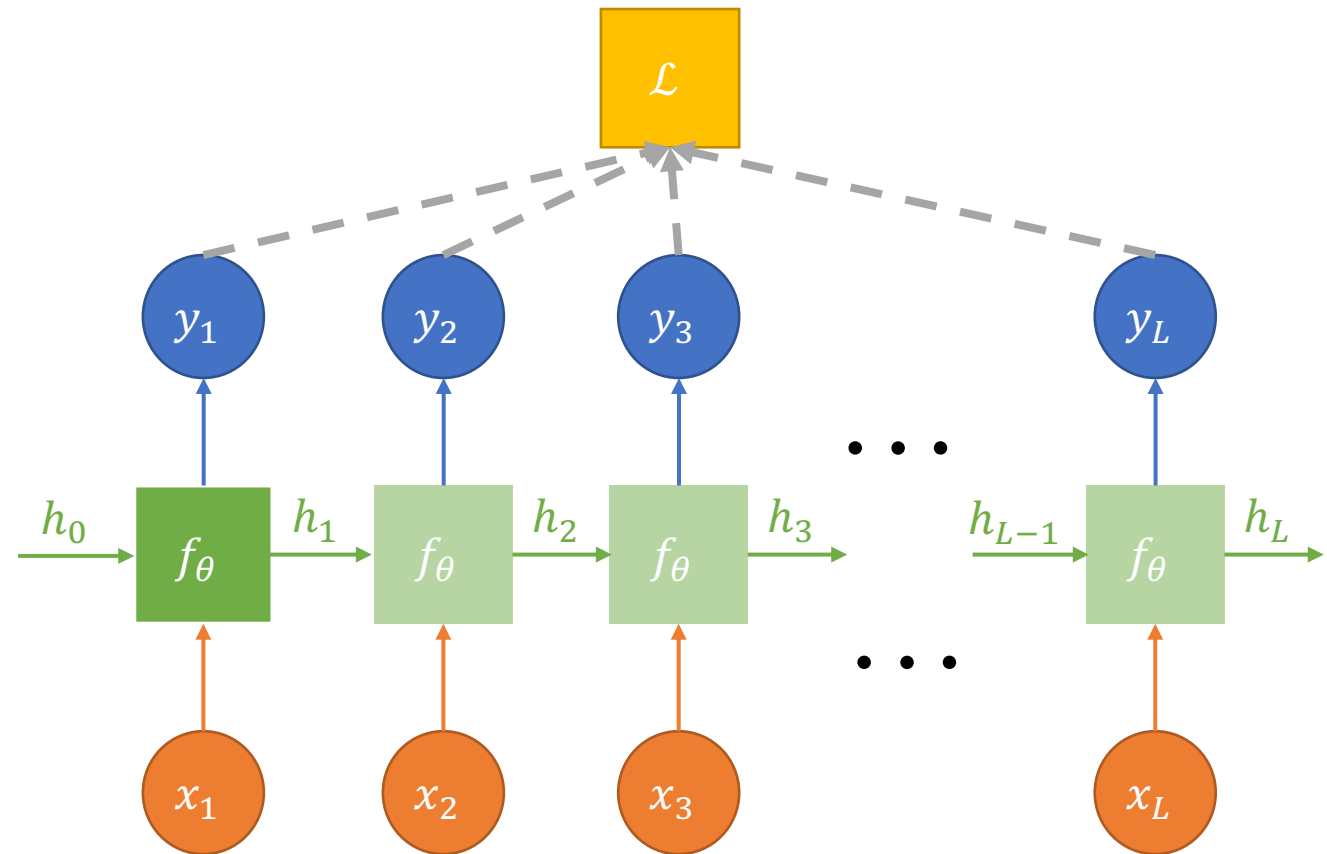
- The RNN module can be written as  $f_{\theta}(h_{\ell-1}, x_{\ell}) = (h_{\ell}, y_{\ell})$ 
  - $h_{\ell} = \tanh(W_h h_{\ell-1} + W_x x_{\ell} + b_h)$
  - $y_{\ell} = W_y h_{\ell} + b_y = W_y \tanh(W_h h_{\ell-1} + W_x x_{\ell}) + b_y$
  - The parameters of the model are the weights and biases  
 $\theta = (W_h, W_x, W_y, b_h, b_y)$





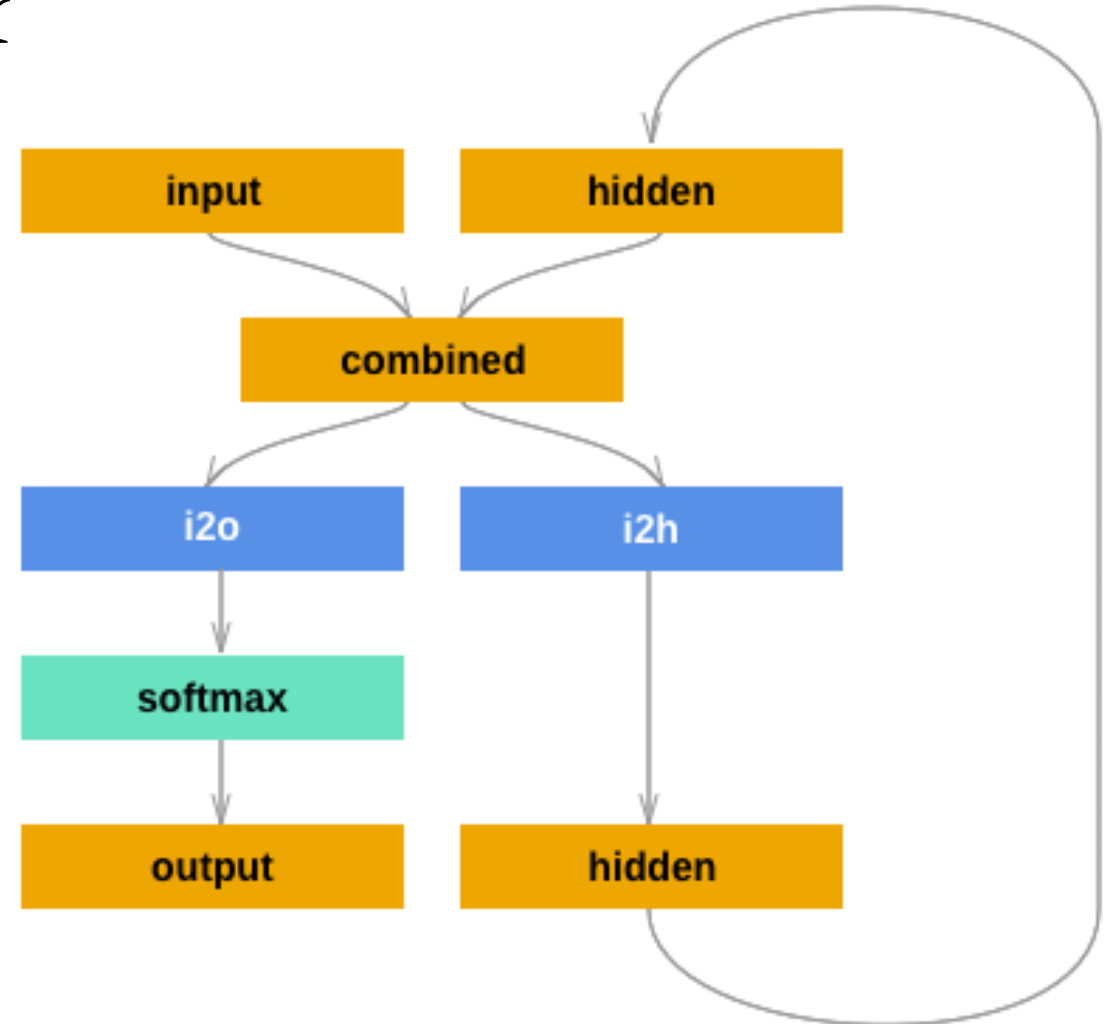
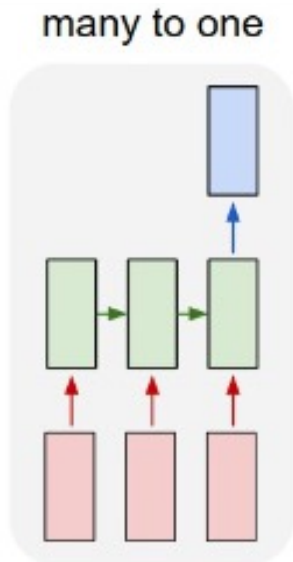
# Training RNNs can simply use backpropagation where the whole chain is “backpropagated”

- Backpropagate gradients of the **loss function** through the computation graph
- PyTorch’s **dynamic** computation graph enables backprop for any length sequence
- Each RNN model evaluation  $f_{\theta}$  is logically the same model so the gradients **accumulate**



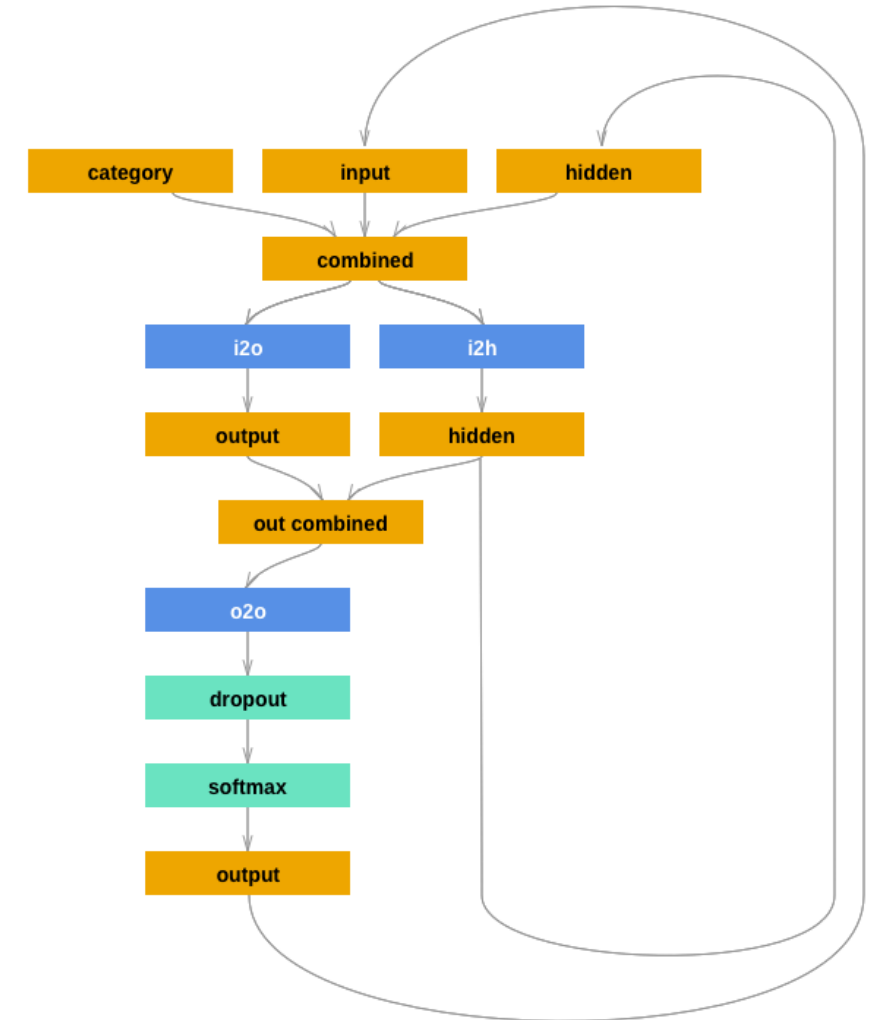
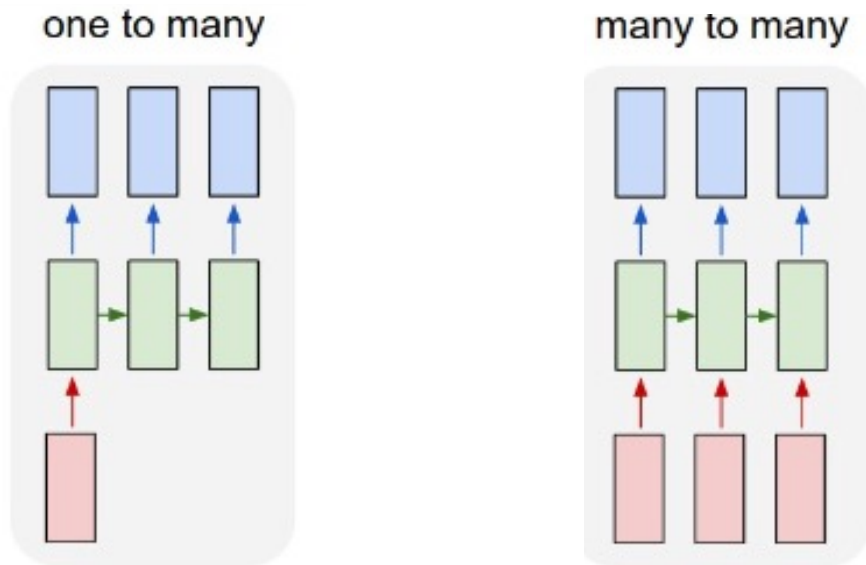
# Demo of simple RNNs for sequence classification

- Task is many to one.
- Architecture is simple.



# Demo of simple RNNs for sequence generation

- Task is conceptually one to many but can be implemented as many to many.
- Architecture is a little more complex.



However, vanilla RNNs suffer from vanishing and exploding gradients that result in only learning short-term dependencies

- Vanishing or exploding gradient are caused by recursive definition of hidden state
- For simplicity, let's assume that  $w_x = 0$  so that we see the core issue.
- The last prediction is as follows:

$$\begin{aligned}\hat{y}_L &= w_y h_L + w_x x_L = w_y h_L = w_y (w_h h_{L-1} + w_x x_{L-1}) \\ &= w_y w_h h_{L-1} = w_y w_h^2 h_{L-2} = \dots = w_y w_h^L h_0\end{aligned}$$

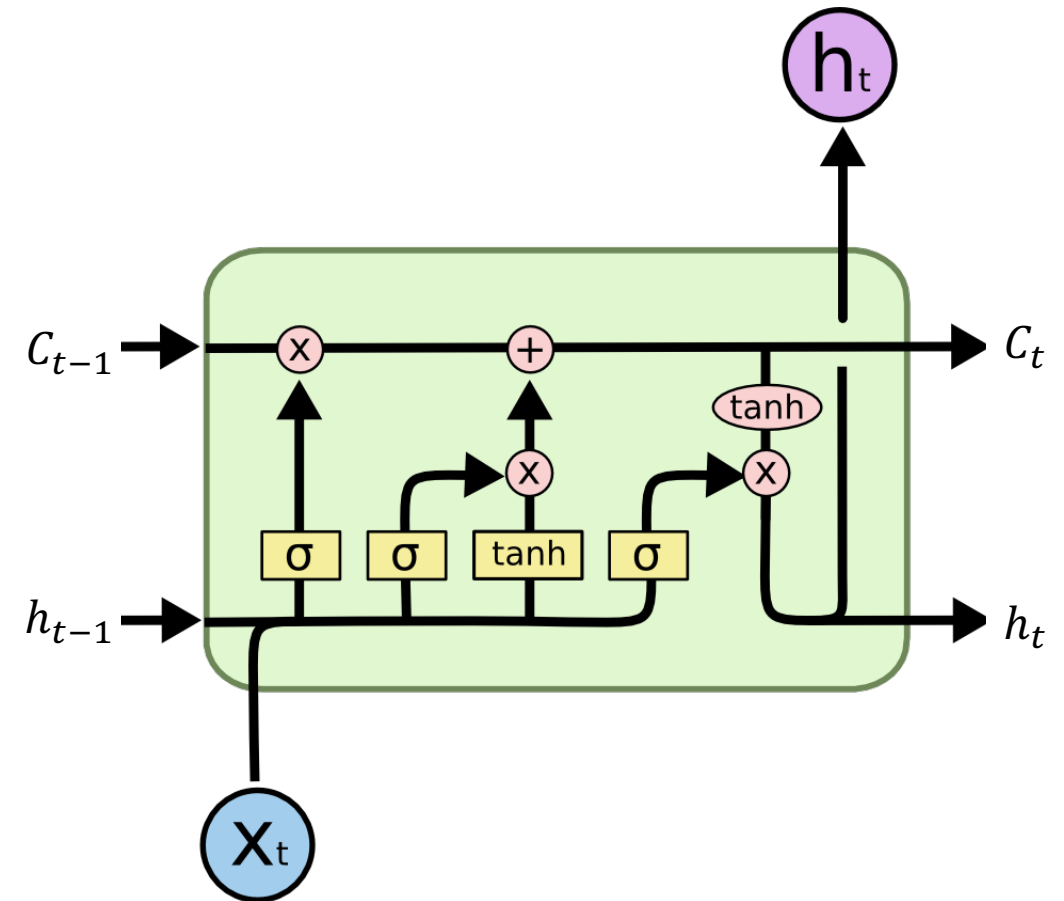
- The gradient of MSE loss for the last term is:

$$\frac{d}{dw_y} \ell(y, \hat{y}_L) = \frac{d}{dw_y} \|y - \hat{y}_L\|_2^2 = 2(y - \hat{y}_L) \frac{d\hat{y}_L}{dw_y} = 2(y - \hat{y}_L) w_h^L h_0$$

- If  $w_h > 1.0$ , then the gradient **exponentially increases** w.r.t. sequence length  $L$ .
- If  $w_h < 1.0$ , then the gradient **exponentially decreases** w.r.t. sequence length  $L$ .
- See demo on simple RNN.

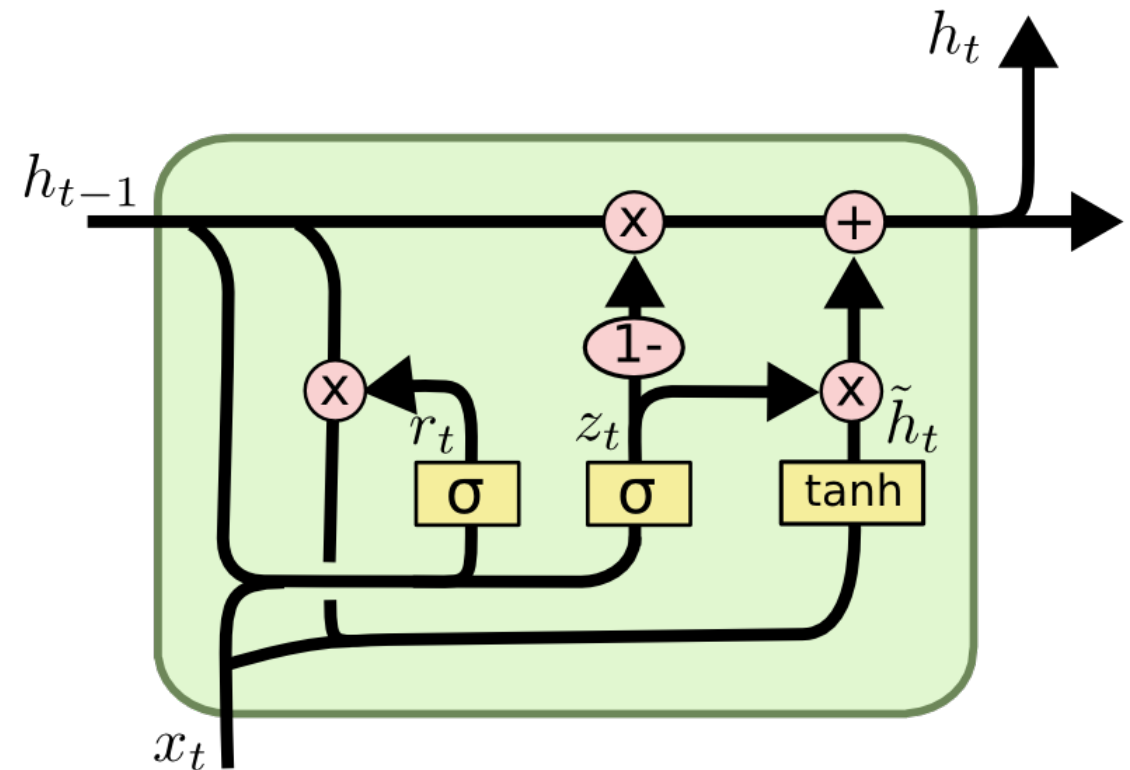
# Long Short-Term Memory (LSTM) units alleviate vanishing gradient problem and enable learning of long-term dependencies

- $h'_{t-1} = [h_{t-1}, x_t]$  (concatenate)
- $\tilde{C}_t = \tanh(W_C h'_{t-1} + b_C)$   
(new cell state information)
- $f_t = \sigma(W_f h'_{t-1} + b_f)$  (forget gate)
- $i_t = \sigma(W_i h'_{t-1} + b_i)$  (input gate)
- $C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$   
(update cell state)
- $o_t = \sigma(W_o h'_{t-1} + b_o)$  (output gate)
- $h_t = o_t \odot \tanh(C_t)$

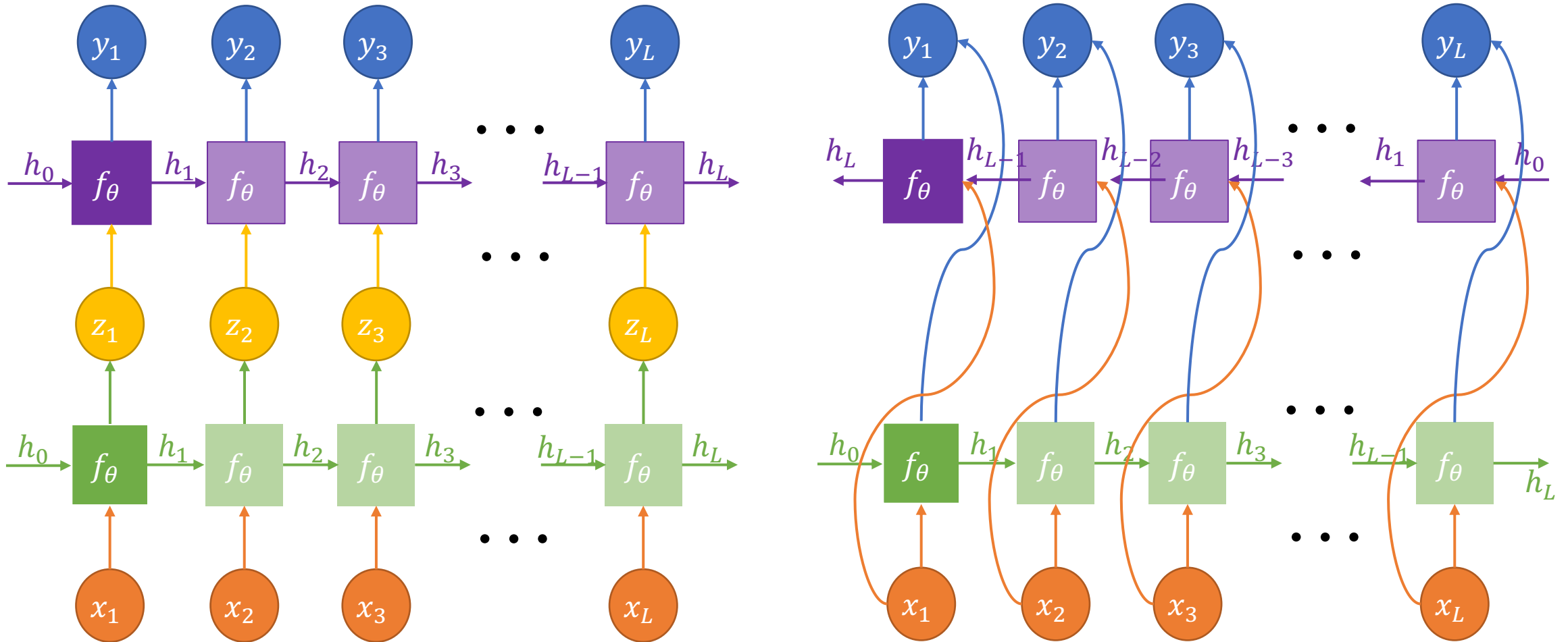


# Gated Recurrent Units (GRU) simplify the LSTM structure and seem to have better performance

- $h'_{t-1} = [h_{t-1}, x_t]$  (concatenate)
- $z_t = \sigma(W_z h'_{t-1})$  (forget/input gate)
- $r_t = \sigma(W_r h'_{t-1})$  (hidden gate)
- $\tilde{h}_t = \tanh(W[r_t \odot h_{t-1}, x_t])$   
(new hidden information)
- $h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$   
(update hidden)



# RNNs can be stacked into deep RNNs and even bidirectional RNNs



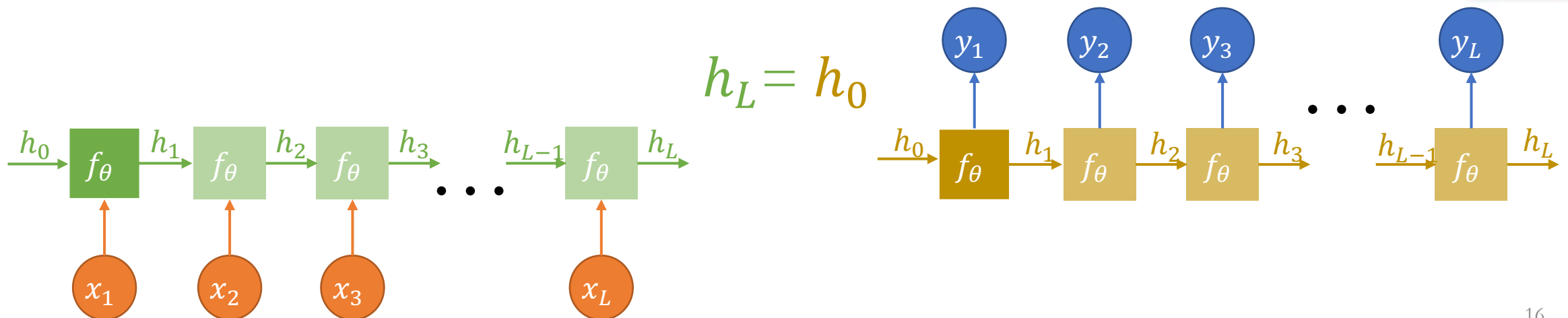
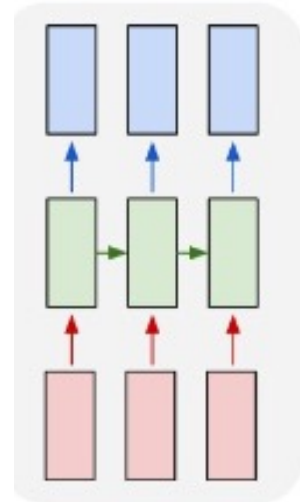
Deep RNN

Bidirectional RNN

# Standard RNNs struggle for sequence-to-sequence tasks because of limited hidden state capacity

- Example: Translation between French and English
- Could we use a one-to-one input/output RNN?
  - Problem: Input sequence could have different length.
  - Problem: The order of words is not the same in French and English.
- More common to use autoencoder structure with 2 RNNs.
  - Problem: Challenging to encode entire sentence in hidden state.

many to many





**Attention** is a model architecture that enables the decoder to efficiently use all encoder outputs

- Attention overcomes some of the challenges of RNN-based translation
- Attention allows long-range dependencies and avoids a completely sequential view of the input and output
- Details will be in next lecture on attention and transformers.

