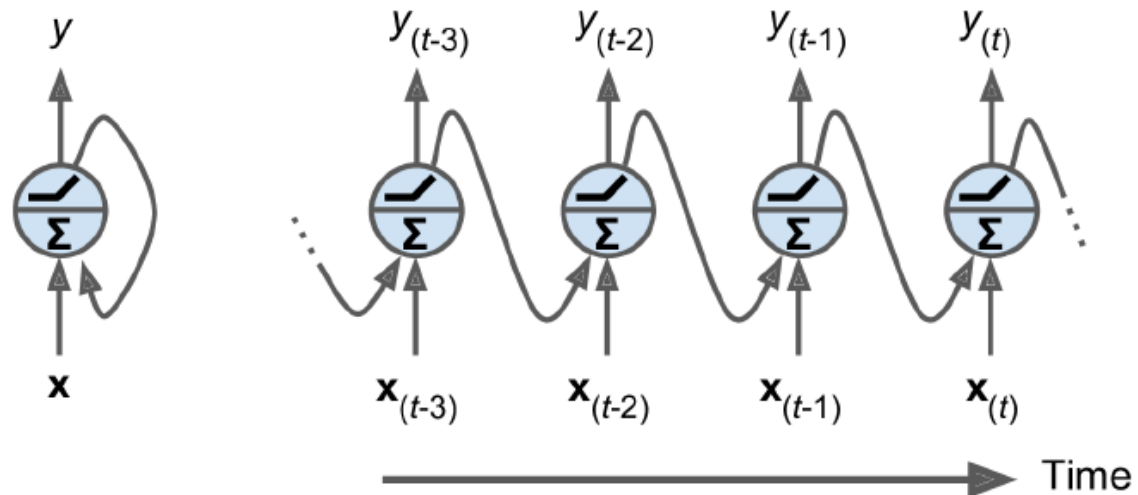# *Recurrent Neural Networks (RNN)*

Saeed Sharifian

# Recurrent Neural Networks

▶ The idea behind Recurrent neural networks (RNN) is to make use of sequential data.
  - Until here, we assume that all inputs (and outputs) are independent of each other.
  - It is a bad idea for many tasks, e.g., predicting the next word in a sentence (it's better to know which words came before it).

▶ They can analyze time series data and predict the future.

▶ They can work on sequences of arbitrary lengths, rather than on fixed-sized inputs.

▶ Neurons in an RNN have connections pointing backward.

▶ RNNs have memory, which captures information about what has been calculated so far.

# Recurrent Neural Networks

We can process a sequence of vectors **x** by applying a **recurrence formula** at every time step:
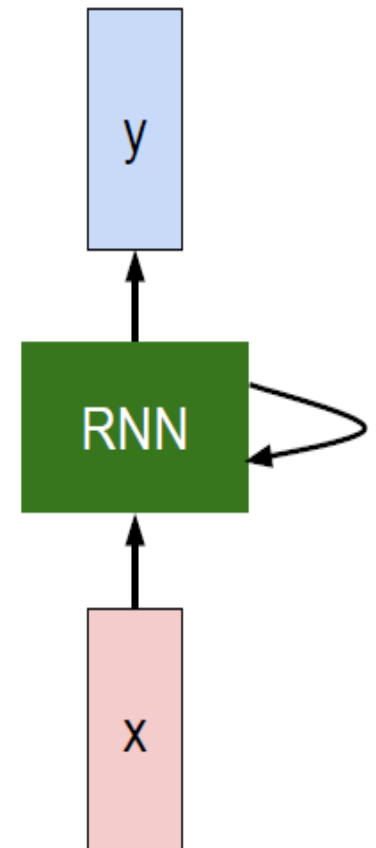
$$h_t = f_W(h_{t-1}, x_t)$$

new state

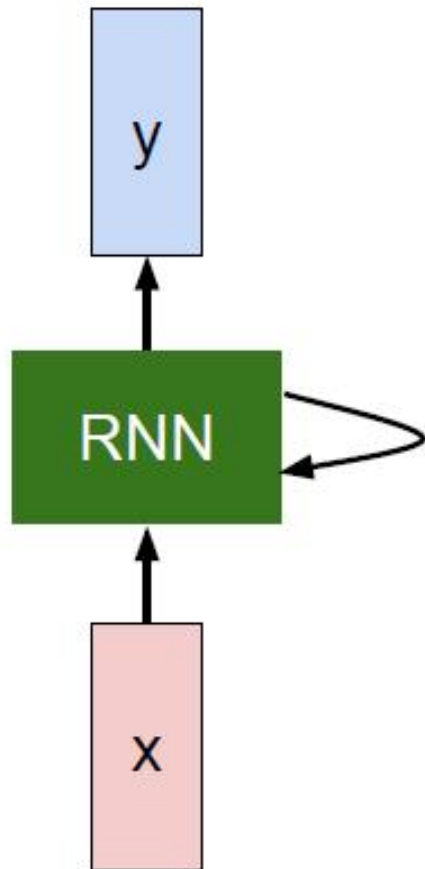some function with parameters W

old state

input vector at some time step

Notice: the same function and the same set of parameters are used at every time step.

# Recurrent Neural Networks

The state consists of a single *"hidden"* vector **h**:
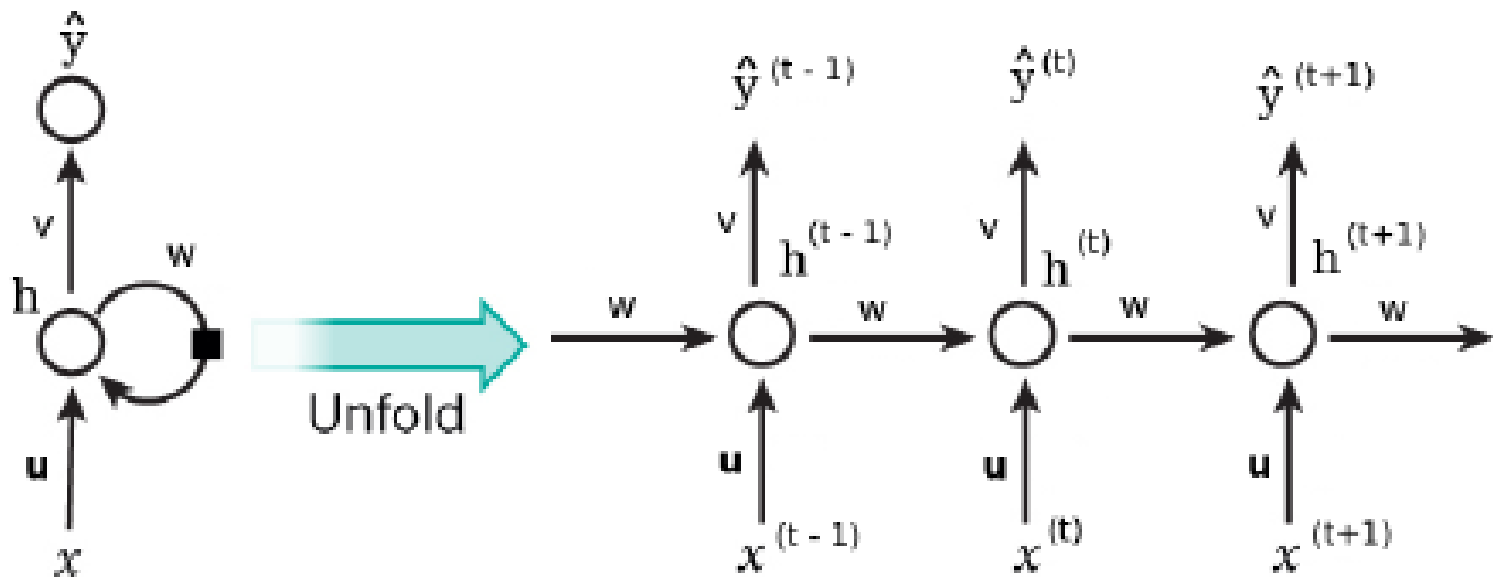


$$h_t = f_W(h_{t-1}, x_t)$$
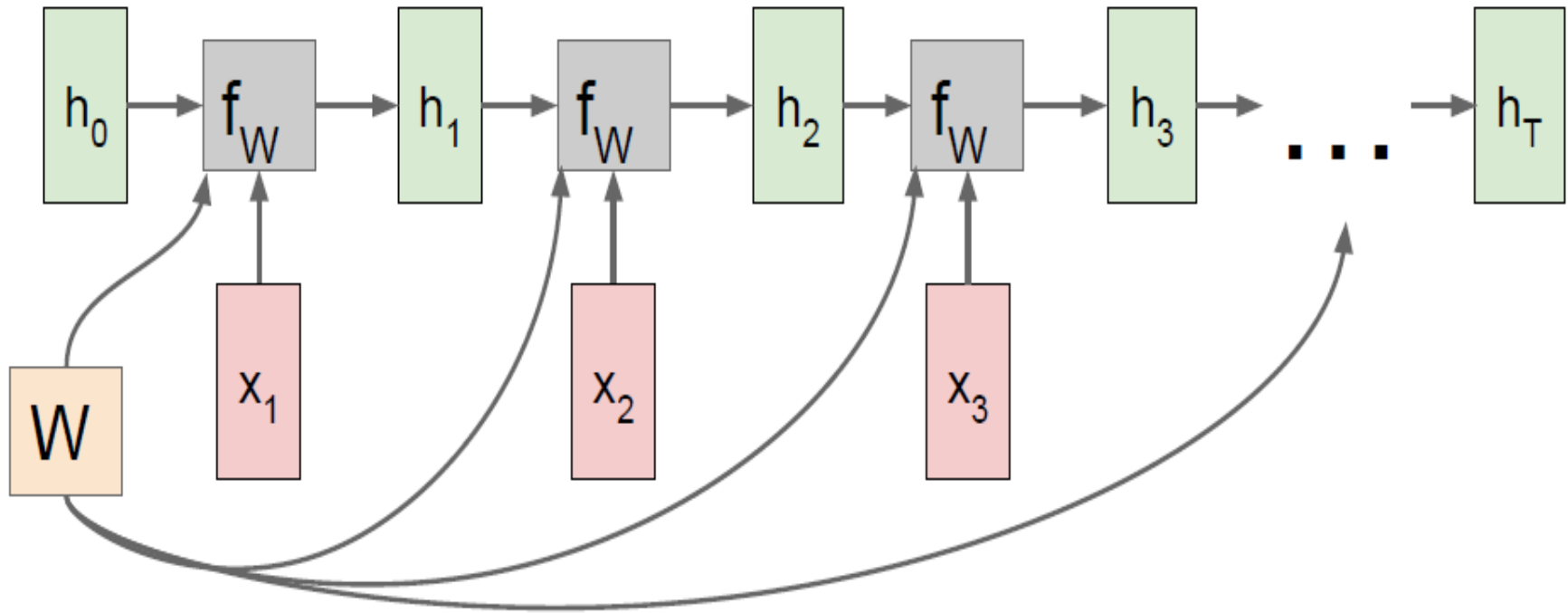
$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$y_t = W_{hy} h_t$$

# Recurrent Neural Networks

▶ Unfolding the network: represent a network against the time axis.
  • We write out the network for the complete sequence.

▶ For example, if the sequence we care about is a sentence of three words, the network would be unfolded into a 3-layer neural network.
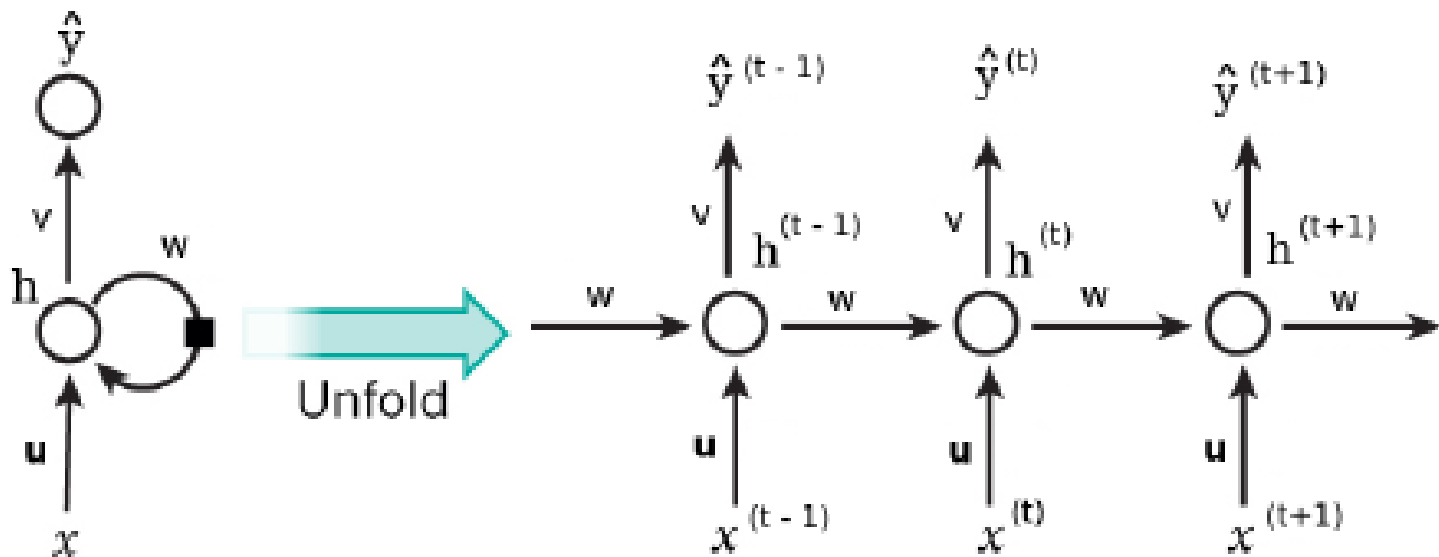  • One layer for each word.

# Recurrent Neural Networks
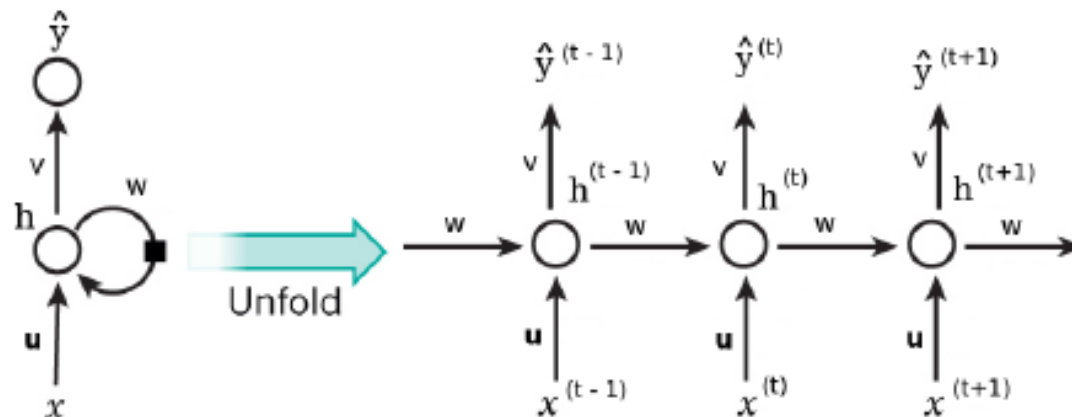


Re-use the same weight matrix at every time-step

# Recurrent Neural Networks

- $h^{(t)} = f(u^T x^{(t)} + w h^{(t-1)})$, where $f$ is an activation function, e.g., `tanh` or `ReLU`.

- $\hat{y}^{(t)} = g(v h^{(t)})$, where $g$ can be the `softmax` function.

- $cost(y^{(t)}, \hat{y}^{(t)}) = cross\_entropy(y^{(t)}, \hat{y}^{(t)}) = -\sum y^{(t)} \log \hat{y}^{(t)}$

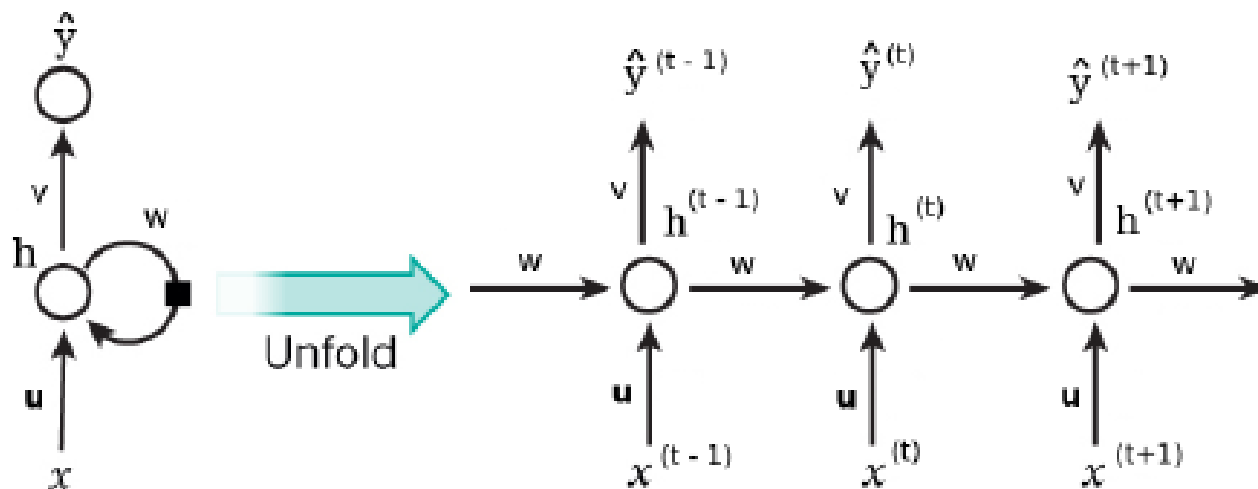- $y^{(t)}$ is the correct word at time step $t$, and $\hat{y}^{(t)}$ is the prediction.

# Recurrent Neurons - Weights

- Each recurrent neuron has three sets of weights: $\mathbf{u}$, $\mathbf{w}$, and $\mathbf{v}$.

- $\mathbf{u}$: the weights for the inputs $\mathbf{x}^{(t)}$.

- $\mathbf{x}^{(t)}$: is the input at time step $t$.

- For example, $\mathbf{x}^{(1)}$ could be a one-hot vector corresponding to the first word of a sentence.

- $\mathbf{w}$: the weights for the hidden state of the previous time step $\mathbf{h}^{(t-1)}$.

- $\mathbf{h}^{(t)}$: is the hidden state (memory) at time step $t$.
  - $\mathbf{h}^{(t)} = \tanh(\mathbf{u}^\mathsf{T}\mathbf{x}^{(t)} + \mathbf{w}\mathbf{h}^{(t-1)})$
  - $\mathbf{h}^{(0)}$ is the initial hidden state.

# Recurrent Neurons - Weights

- v: the weights for the hidden state of the current time step $h^{(t)}$.

- $\hat{y}^{(t)}$ is the output at step $t$.

- $\hat{y}^{(t)} = \text{softmax}(vh^{(t)})$

- For example, if we wanted to predict the next word in a sentence, it would be a vector of probabilities across our vocabulary.
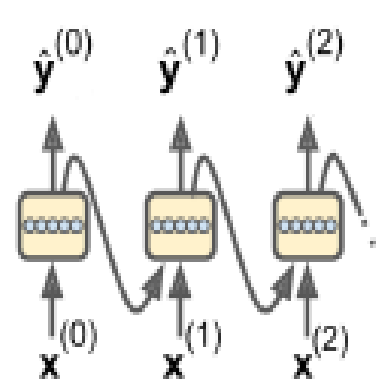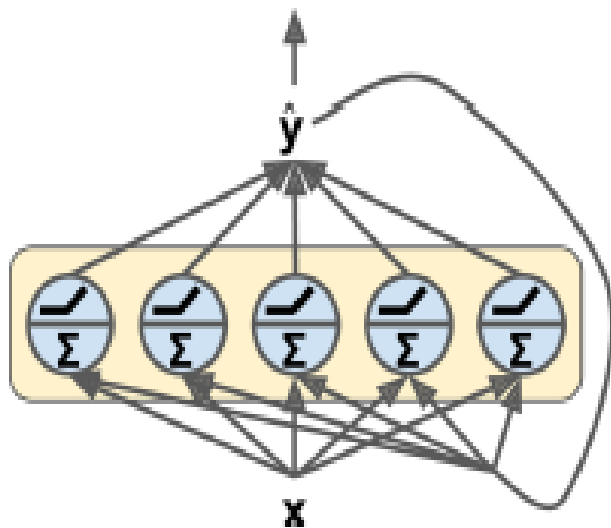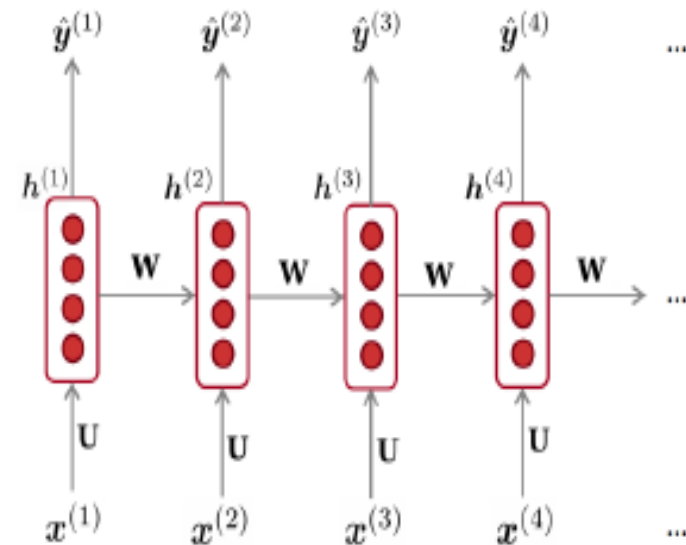
# Layers of Recurrent Neurons

▶ At each time step $t$, every neuron of a layer receives both the input vector $\mathbf{x}^{(t)}$ and the output vector from the previous time step $\mathbf{h}^{(t-1)}$.

$$\mathbf{h}^{(t)} = \tanh(\mathbf{u}^T\mathbf{x}^{(t)} + \mathbf{w}^T\mathbf{h}^{(t-1)})$$

$$\mathbf{y}^{(t)} = \text{sigmoid}(\mathbf{v}^T\mathbf{h}^{(t)})$$

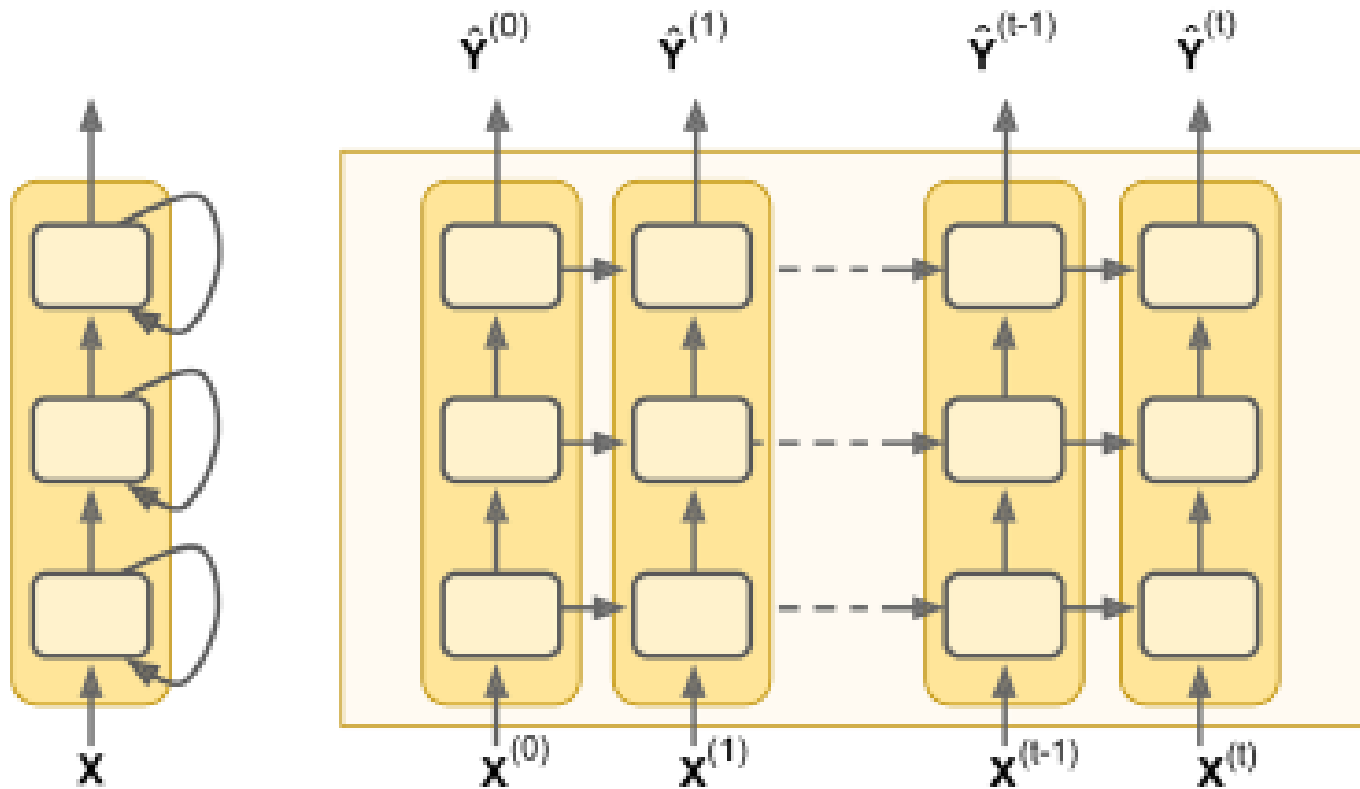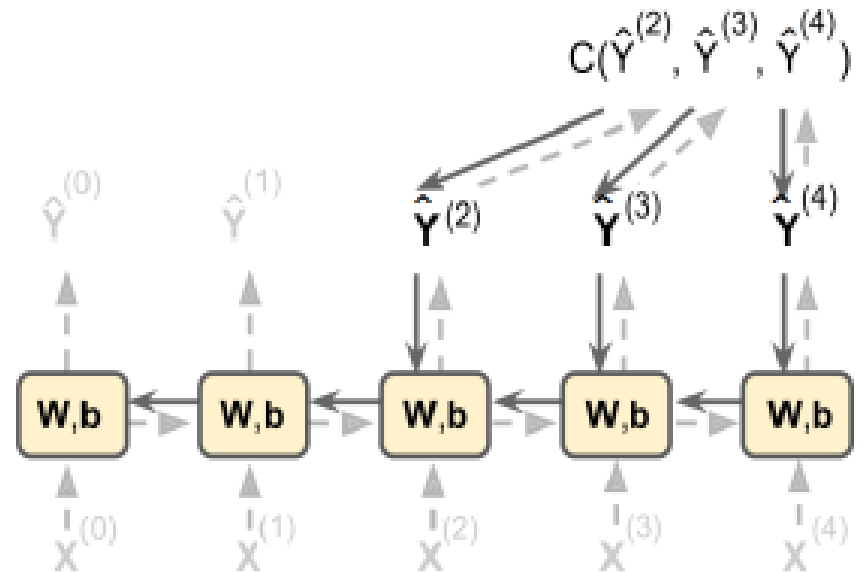# Deep RNN

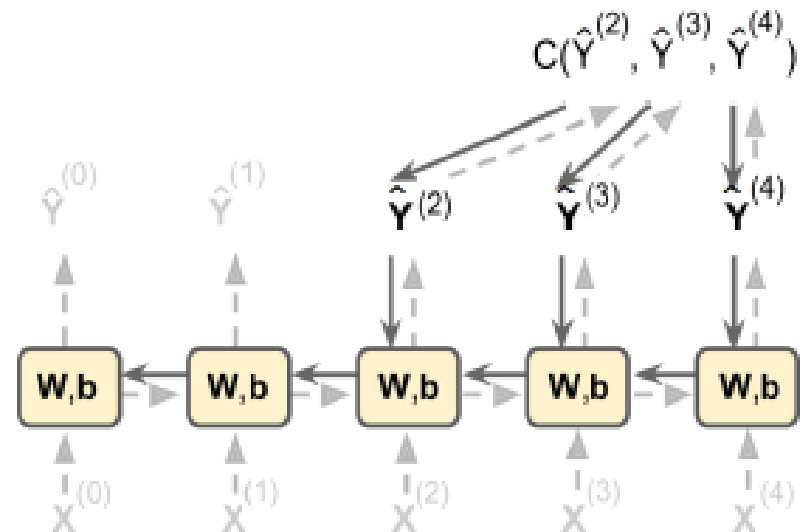▶ Stacking multiple layers of cells gives you a deep RNN.

# Training RNNs

▶ To train an RNN, we should unroll it through time and then simply use regular backpropagation.

▶ This strategy is called backpropagation through time (BPTT).

▶ To train the model using BPTT, we go through the following steps:

▶ 1. Forward pass through the unrolled network (represented by the dashed arrows).

▶ 2. The cost function is $C(\hat{y}^{tmin}, \hat{y}^{tmin+1}, \cdots, \hat{y}^{tmax})$, where tmin and tmax are the first and last output time steps, not counting the ignored outputs.

# Backpropagation Through Time

▶ 3.  Propagate backward the gradients of that cost function through the unrolled network (represented by the solid arrows).

▶ 4.  The model parameters are updated using the gradients computed during BPTT.

▶ The gradients flow backward through all the outputs used by the cost function, not just through the final output.

▶ For example, in the following figure:

- The cost function is computed using the last three outputs, $\hat{\mathbf{y}}^{(2)}$, $\hat{\mathbf{y}}^{(3)}$, and $\hat{\mathbf{y}}^{(4)}$.
- Gradients flow through these three outputs, but not through $\hat{\mathbf{y}}^{(0)}$ and $\hat{\mathbf{y}}^{(1)}$.

# BPTT Step by Step

$$s^{(t)} = u^T x^{(t)} + wh^{(t-1)}$$

$$h^{(t)} = \tanh(s^{(t)})$$

$$z^{(t)} = vh^{(t)}$$

$$\hat{y}^{(t)} = \texttt{softmax}(z^{(t)})$$

$$J^{(t)} = \texttt{cross\_entropy}(y^{(t)}, \hat{y}^{(t)}) = -\sum y^{(t)} \log \hat{y}^{(t)}$$

$$\hat{y}_t = \phi(V h_t)$$

$$h_t = \psi(U x_t + W h_{t-1})$$

▶ We treat the full sequence as one training example.

▶ The total error $E$ is just the sum of the errors at each time step.

▶ E.g., $E = J^{(1)} + J^{(2)} + \cdots + J^{(t)}$

# *RNN Design Patterns*

# Vector-to-Sequence

▶ Vector-to-sequence network: takes a single input at the first time step, and let it output a sequence.

▶ E.g., the input could be an image, and the output could be a caption for that image.

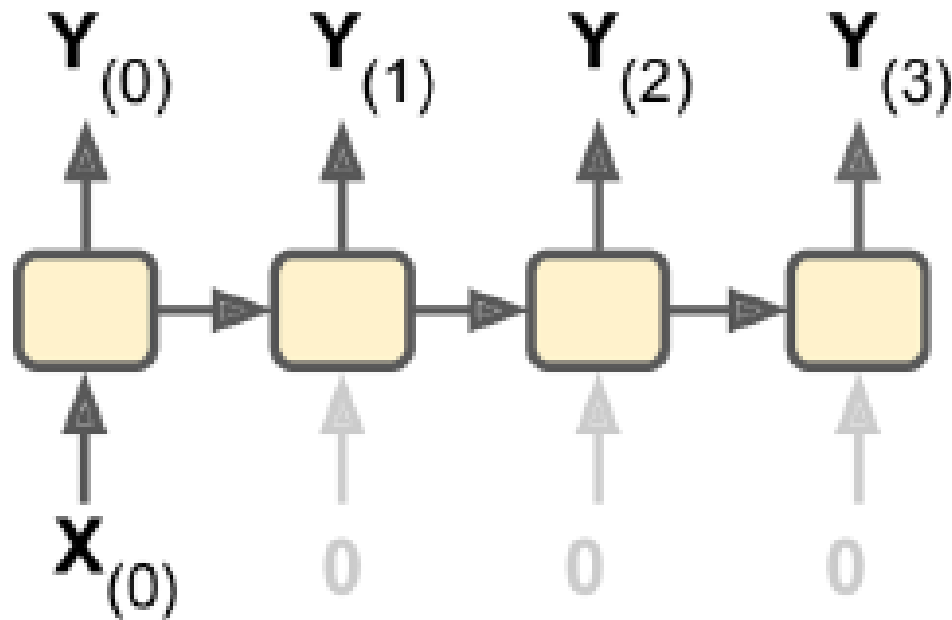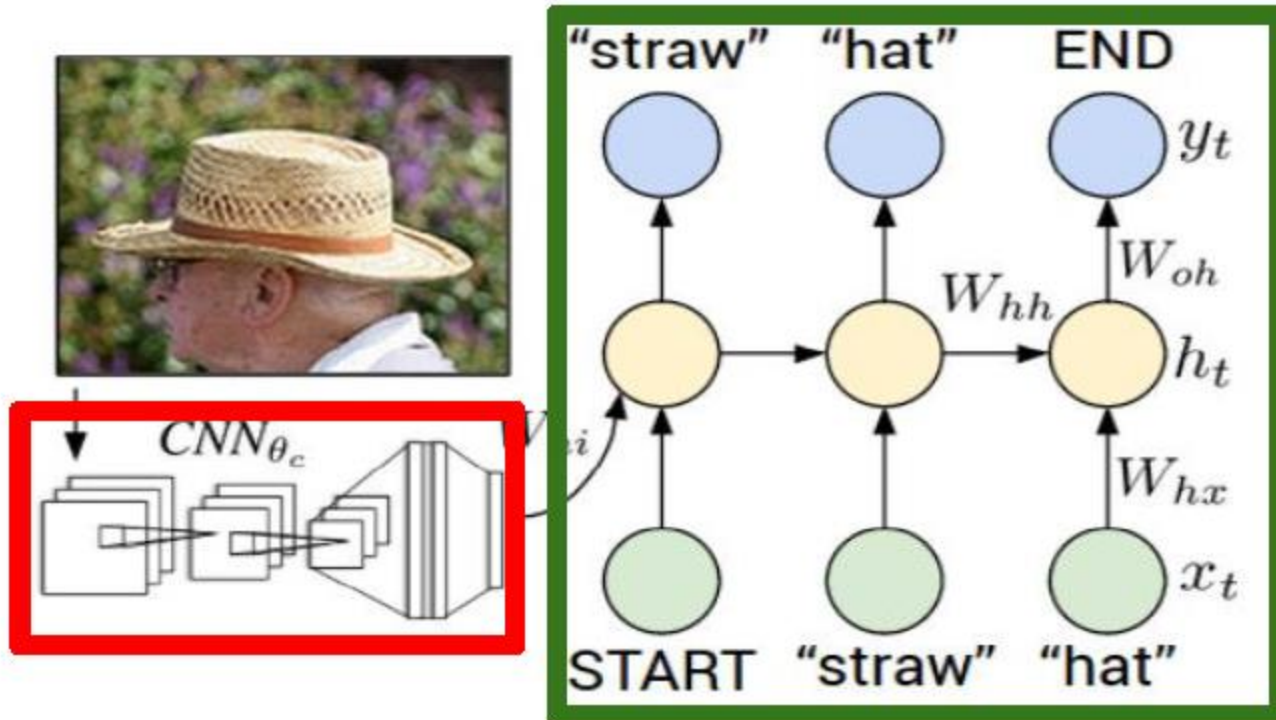$$Y_{(0)} \quad Y_{(1)} \quad Y_{(2)} \quad Y_{(3)}$$

**Image Captioning**
image -> sequence of words

$$X_{(0)} \quad 0 \quad 0 \quad 0$$

# Image Captioning



**Recurrent Neural Network**

**Convolutional Neural Network**

# Image Captioning



test image

**before:**

$h = \tanh(Wxh * x + Whh * h)$

**now:**

$h = \tanh(Wxh * x + Whh * h + \textbf{Wih} * \textbf{v})$

**Wih**

image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096

y0

h0

x0
<START>

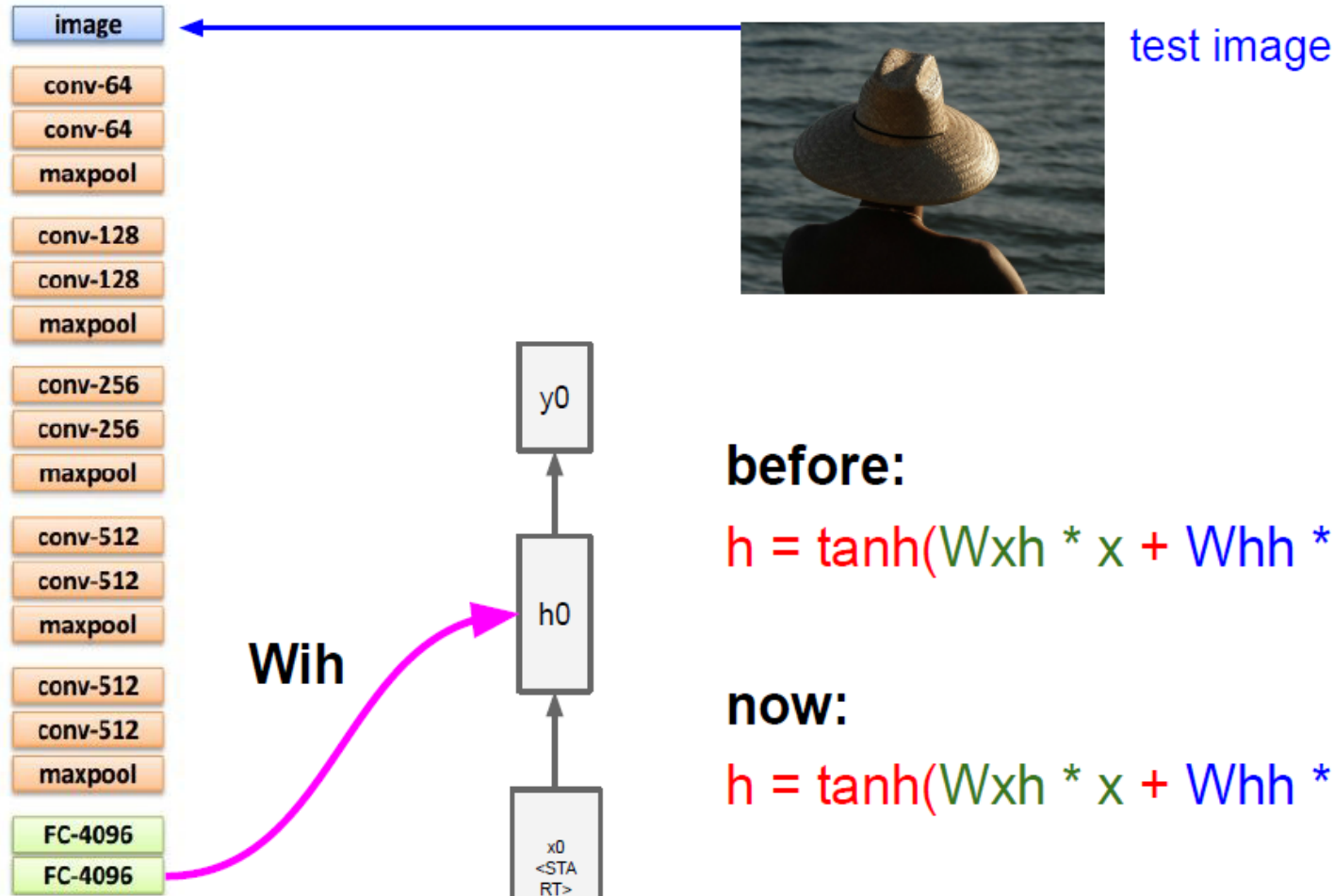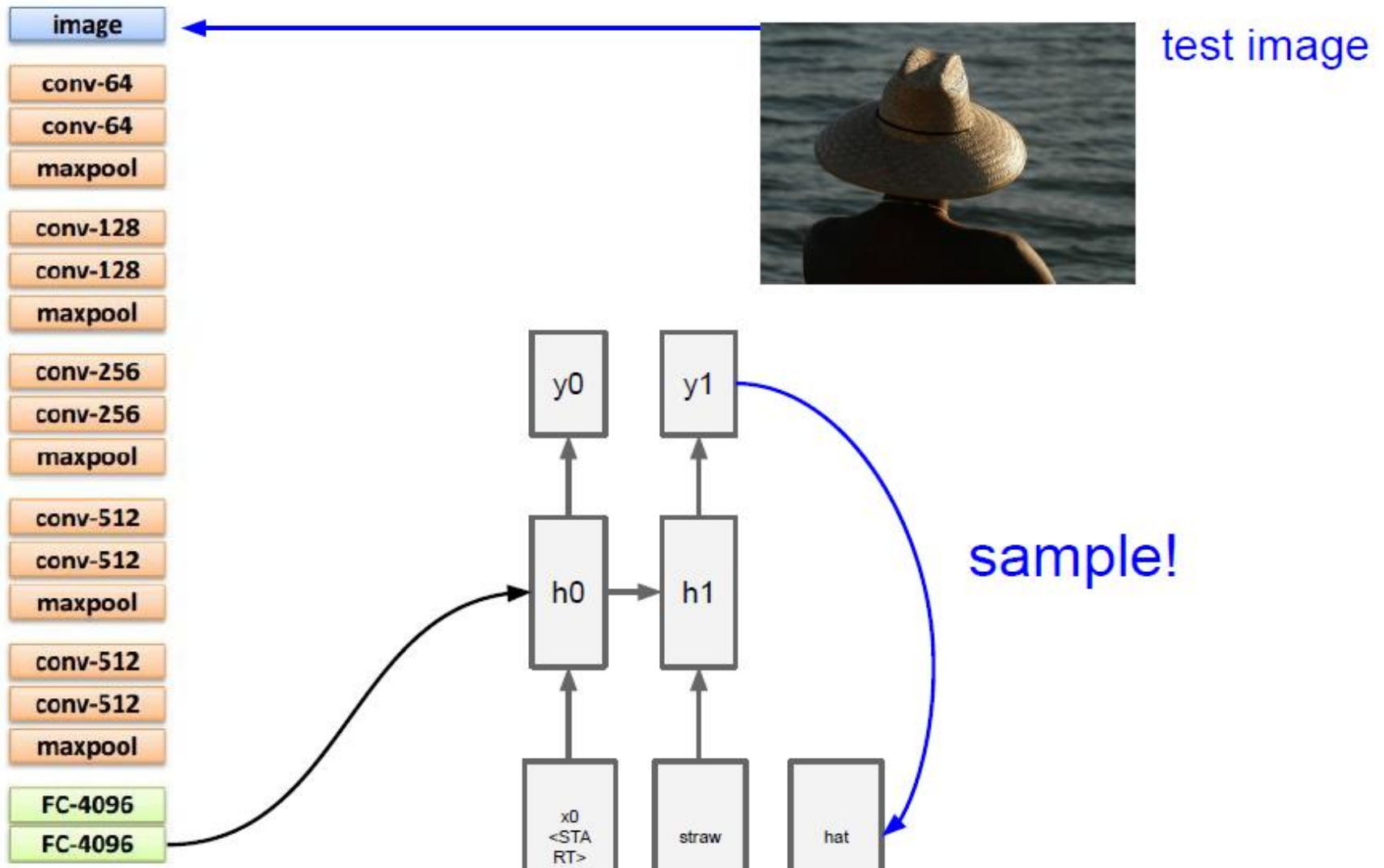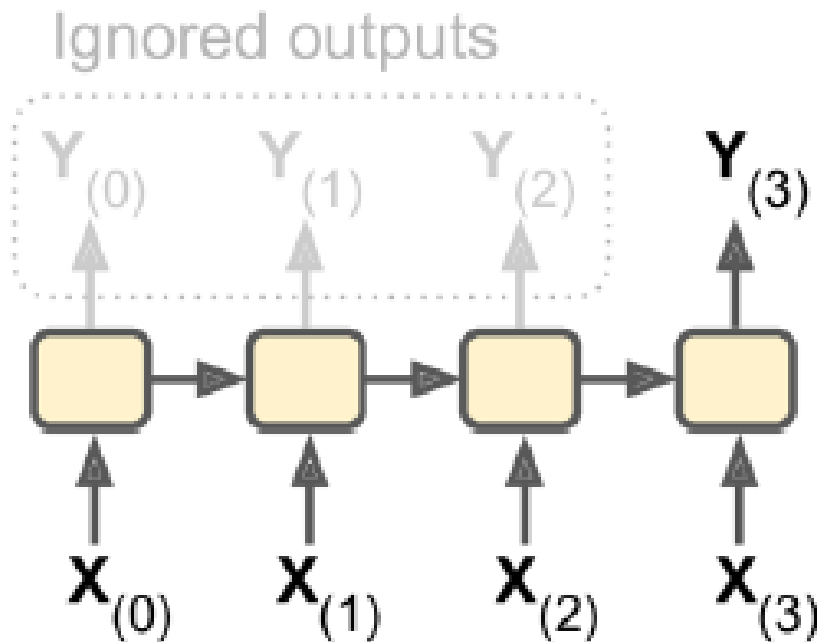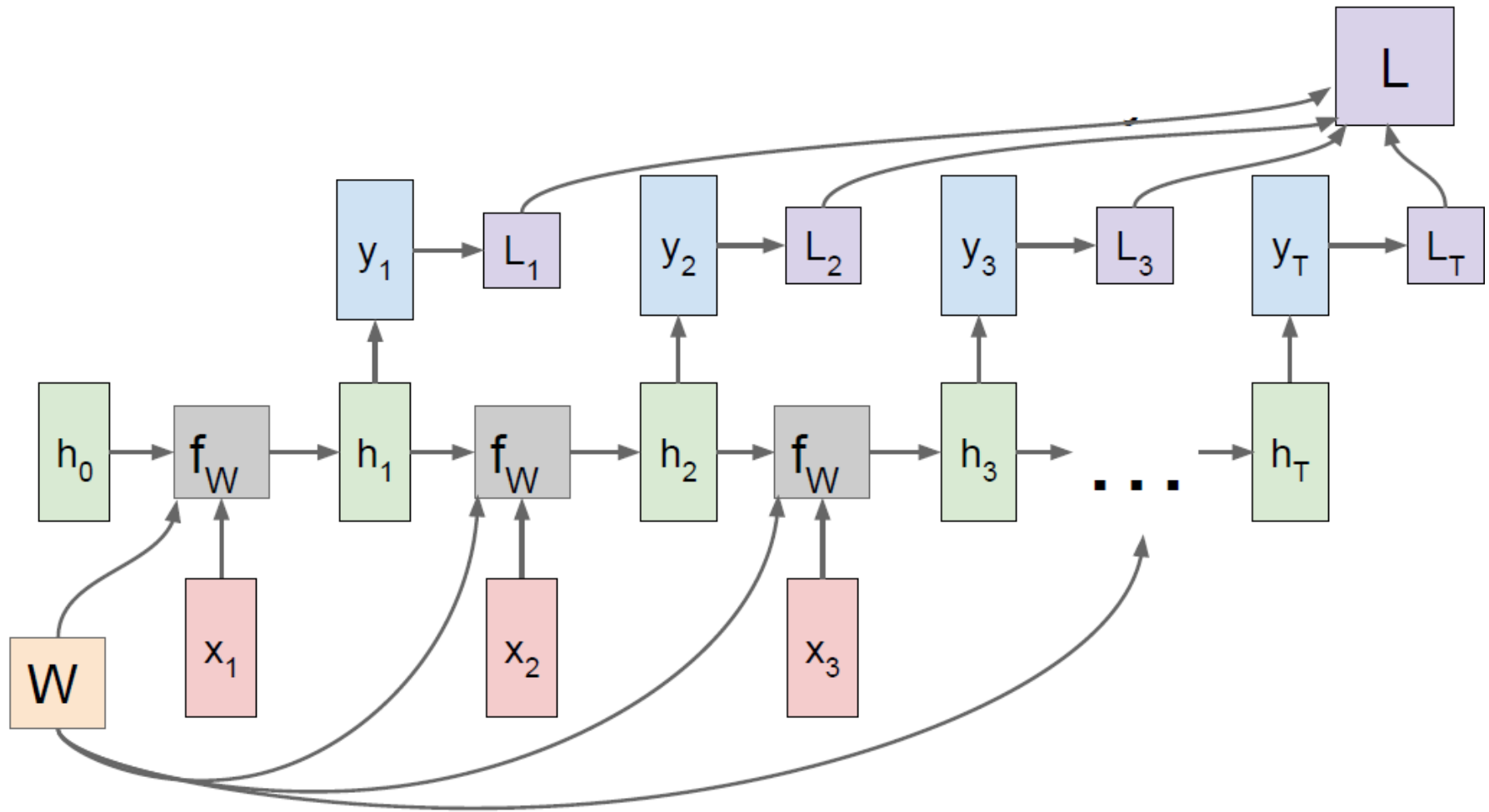# Image Captioning

# Sequence-to-Vector

▶ Sequence-to-vector network: takes a sequence of inputs, and ignore all outputs except for the last one.

▶ E.g., you could feed the network a sequence of words corresponding to a movie review, and the network would output a sentiment score.

Ignored outputs

$Y_{(0)}$    $Y_{(1)}$    $Y_{(2)}$    $Y_{(3)}$
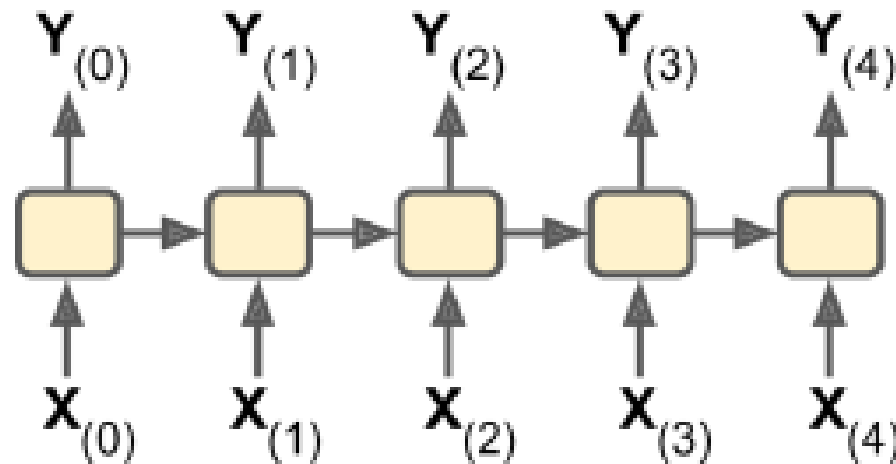
$X_{(0)}$    $X_{(1)}$    $X_{(2)}$    $X_{(3)}$

**Sentiment Classification**
sequence of words -> sentiment

# Many to Many

# Sequence-to-Sequence
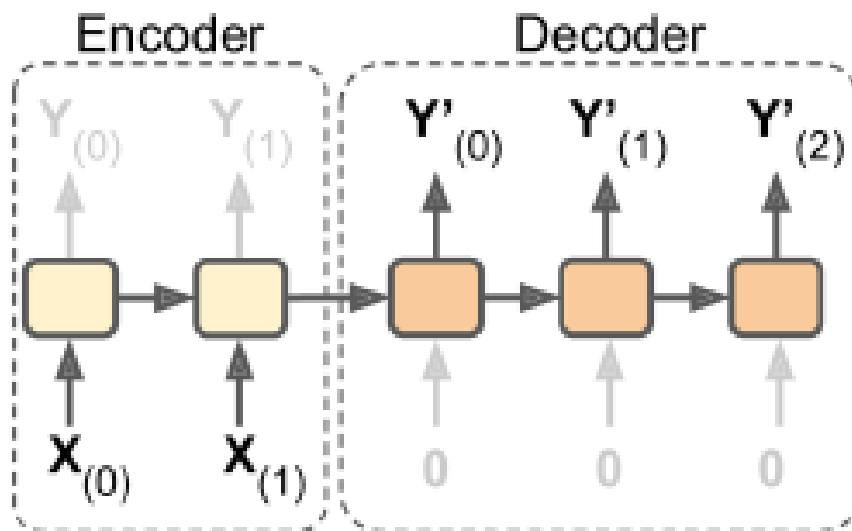
▶ Sequence-to-sequence network: takes a sequence of inputs and produce a sequence of outputs.

▶ Useful for predicting time series such as stock prices: you feed it the prices over the last N days, and it must output the prices shifted by one day into the future.

▶ Here, both input sequences and output sequences have the same length.

$Y_{(0)}$  $Y_{(1)}$  $Y_{(2)}$  $Y_{(3)}$  $Y_{(4)}$

$X_{(0)}$  $X_{(1)}$  $X_{(2)}$  $X_{(3)}$  $X_{(4)}$

**Video classification on frame level**

# Encoder-Decoder

- Encoder-decoder network: a sequence-to-vector network (encoder), followed by a vector-to-sequence network (decoder).

- E.g., translating a sentence from one language to another.

- You would feed the network a sentence in one language, the encoder would convert this sentence into a single vector representation, and then the decoder would decode this vector into a sentence in another language.
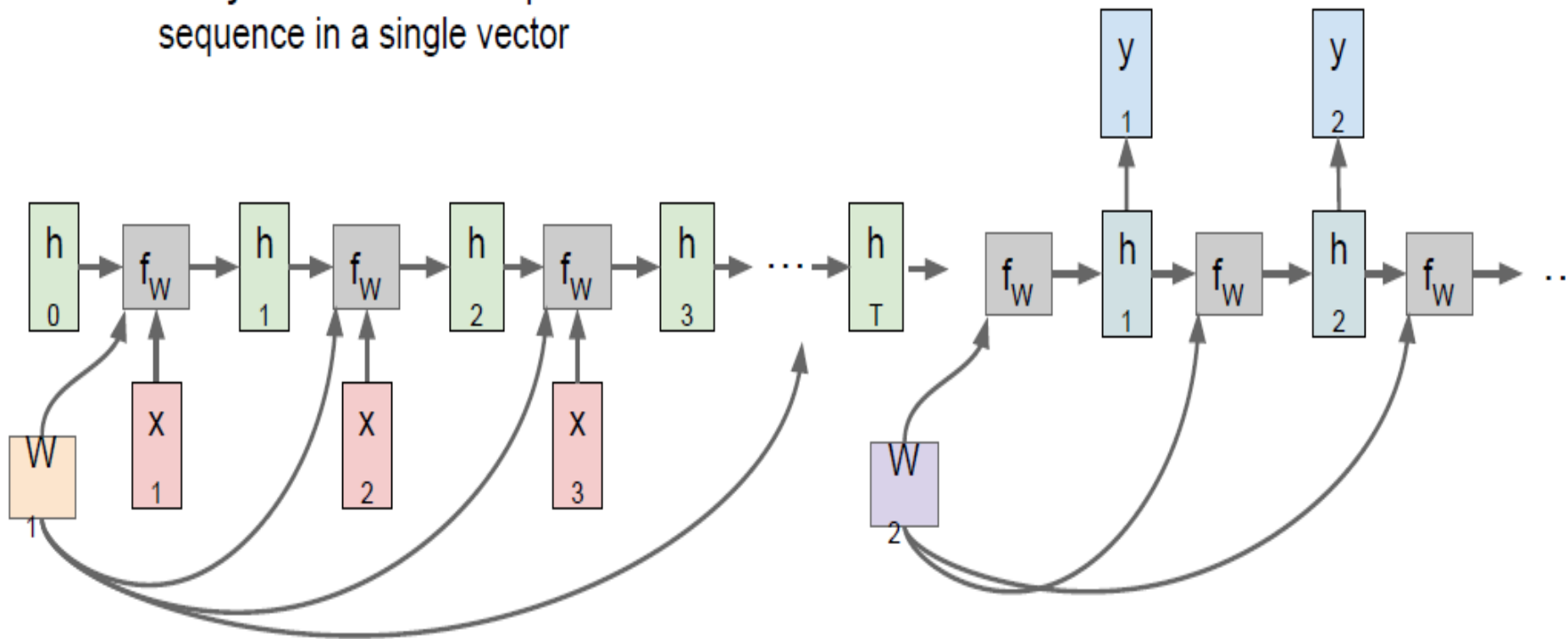


**Machine Translation**
seq of words -> seq of words

# Many-to-one +one-to-many

**Many to one**: Encode input sequence in a single vector

**One to many**: Produce output sequence from single input vector
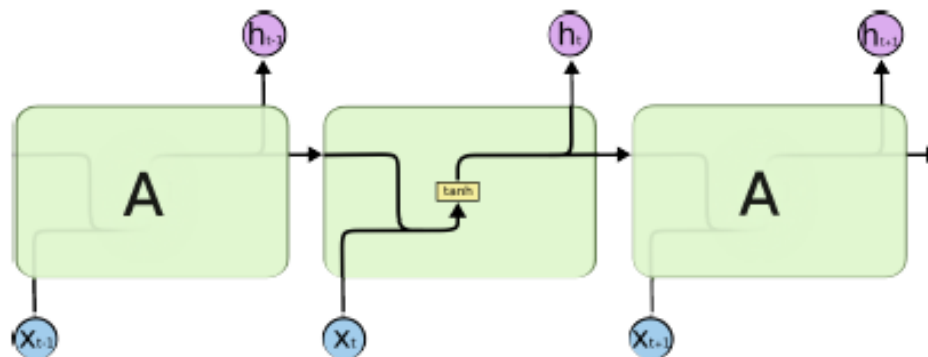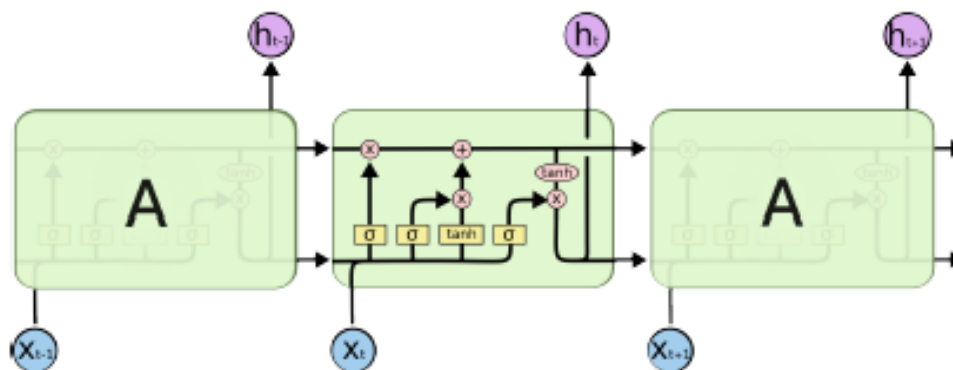
# *long short-term memory (LSTM)*

# RNN Problems

▶ Sometimes we only need to look at recent information to perform the present task.
  - E.g., predicting the next word based on the previous ones.

▶ In such cases, where the gap between the relevant information and the place that it's needed is small, RNNs can learn to use the past information.

▶ But, as that gap grows, RNNs become unable to learn to connect the information.

▶ RNNs may suffer from the vanishing/exploding gradients problem.

▶ To solve these problem, long short-term memory (LSTM) have been introduced.

▶ In LSTM, the network can learn what to store and what to throw away.

# RNN Basic Cell vs. LSTM

▸ Without looking inside the box, the LSTM cell looks exactly like a basic cell.

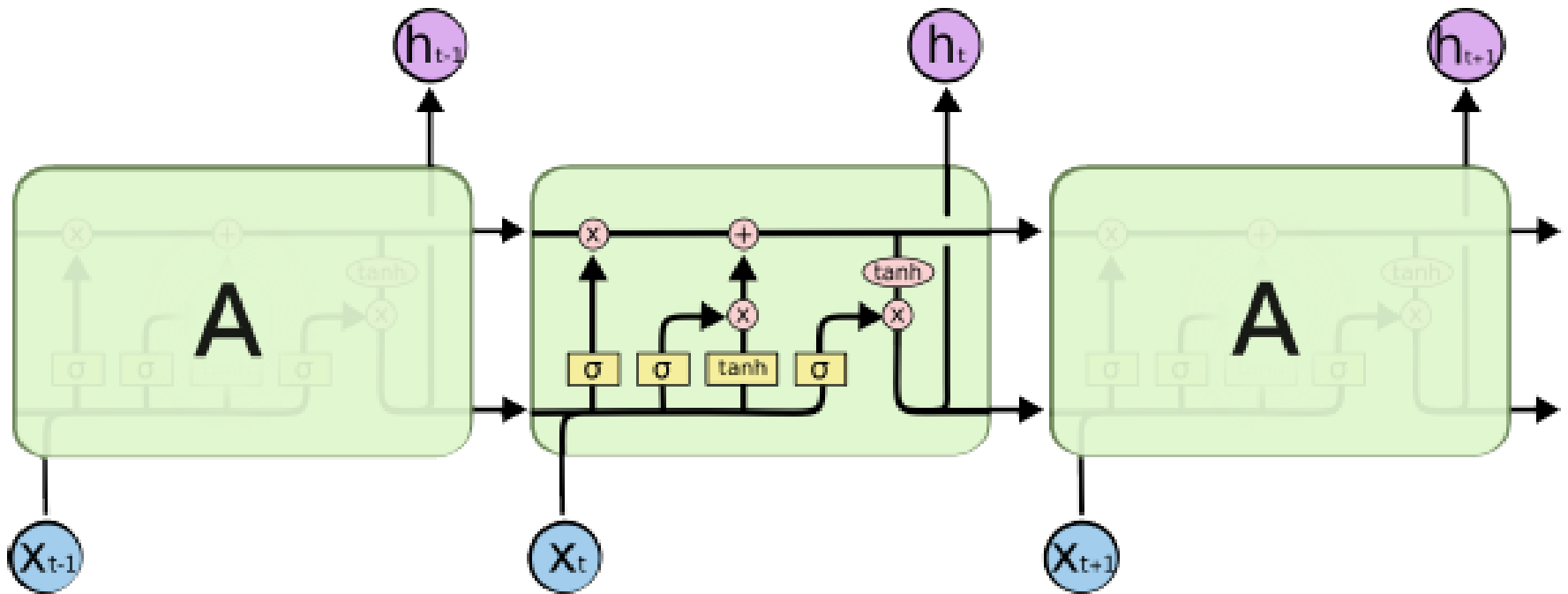▸ The repeating module in a standard RNN contains a single layer.

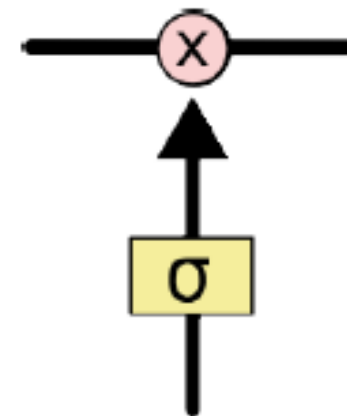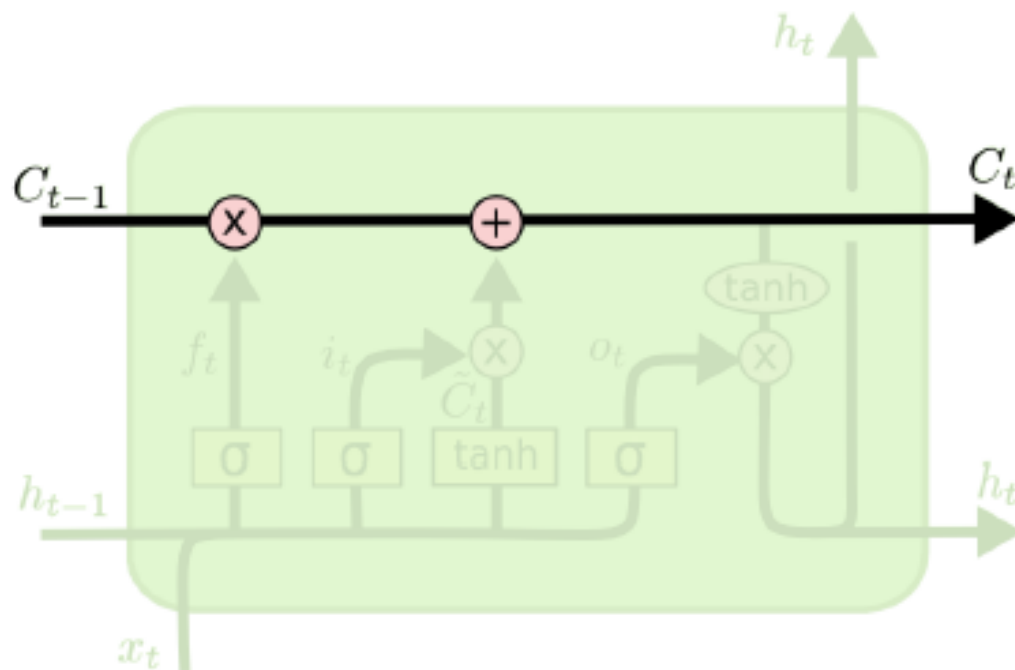▸ The repeating module in an LSTM contains four interacting layers.

# LSTM

▶ In LSTM state is split in two vectors:
  1. $h^{(t)}$ ($h$ stands for hidden): the short-term state
  2. $c^{(t)}$ ($c$ stands for cell): the long-term state

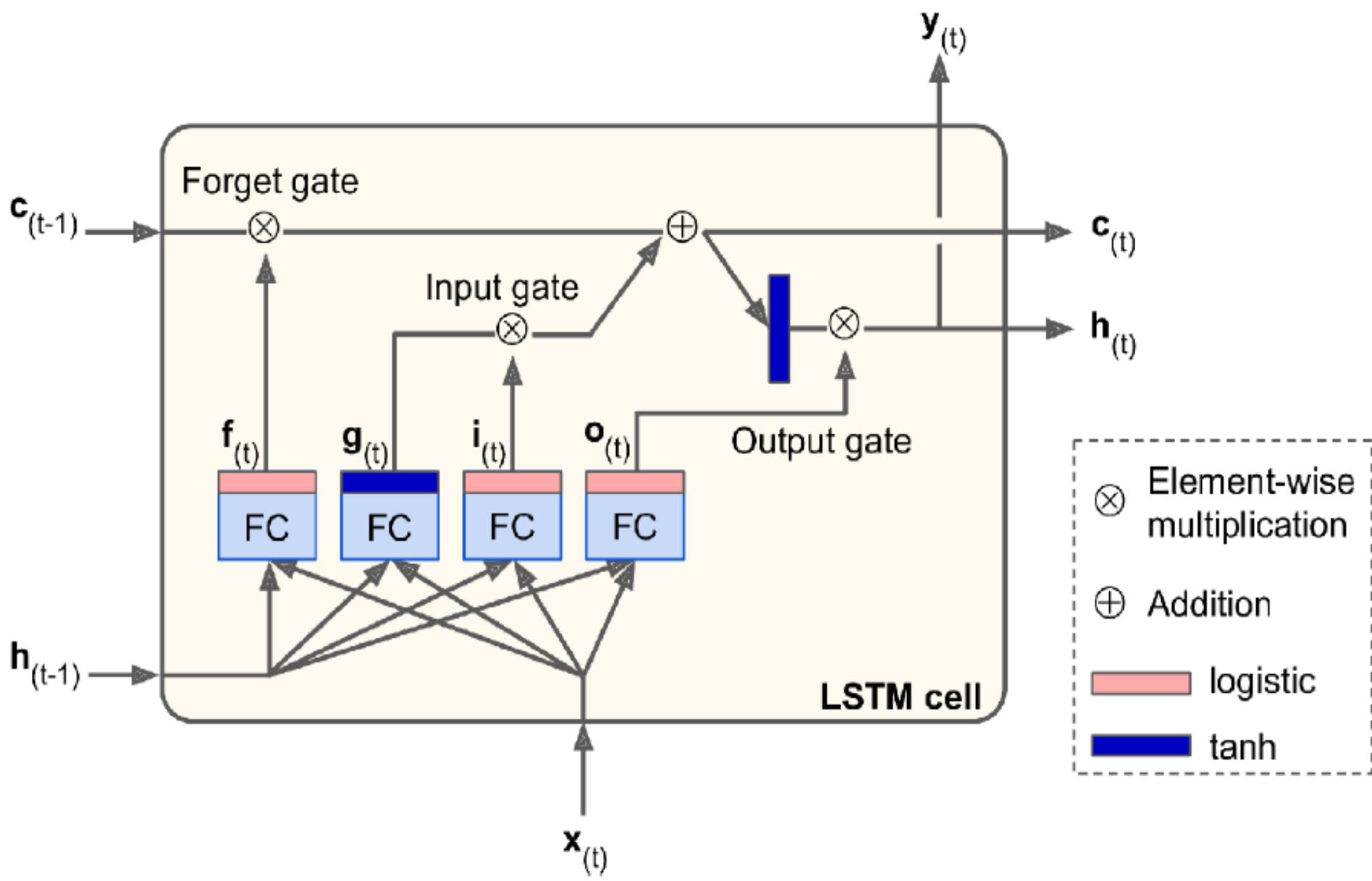# LSTM

▶ The cell state (long-term state), the horizontal line on the top of the diagram.

▶ The LSTM can remove/add information to the cell state, regulated by three gates.
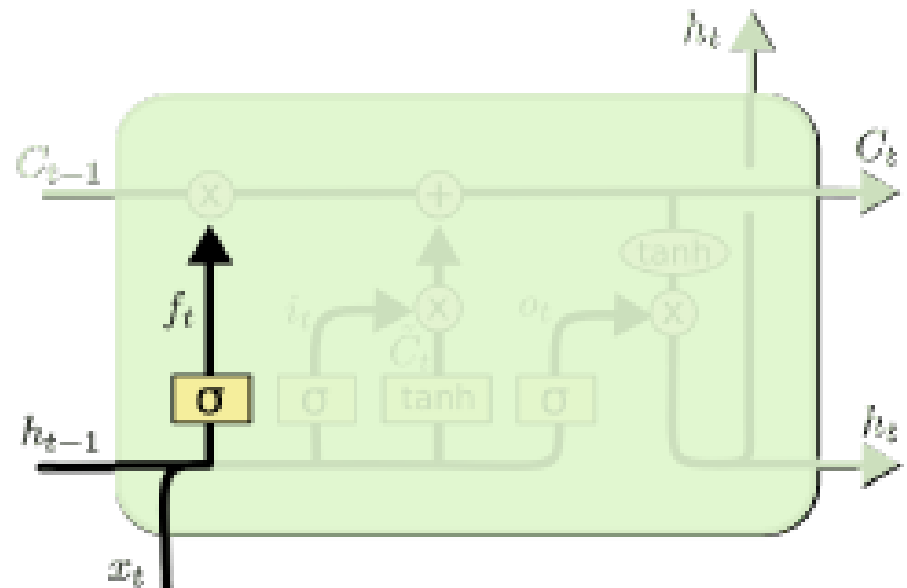
  • Forget gate, input gate and output gate

# LSTM

# Step-by-Step LSTM Walk Through

▸ **Step one**: decides what information we are going to throw away from the cell state.

▸ This decision is made by a sigmoid layer, called the forget gate layer.

▸ It looks at $h^{(t-1)}$ and $x^{(t)}$, and outputs a number between 0 and 1 for each number in the cell state $c^{(t-1)}$.

  • 1 represents completely keep this, and 0 represents completely get rid of this.

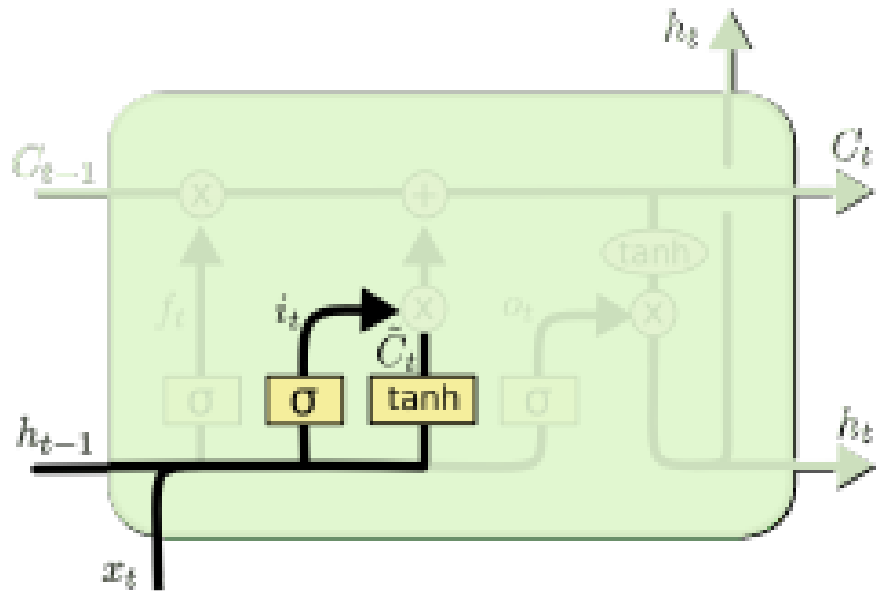$$f^{(t)} = \sigma(u_f^\top x^{(t)} + w_f h^{(t-1)})$$

# Step-by-Step LSTM Walk Through

▶ Second step: decides what new information we are going to store in the cell state. This has two parts:

▶ 1. A sigmoid layer, called the input gate layer, decides which values we will update.

▶ 2. A tanh layer creates a vector of new candidate values that could be added to the state.

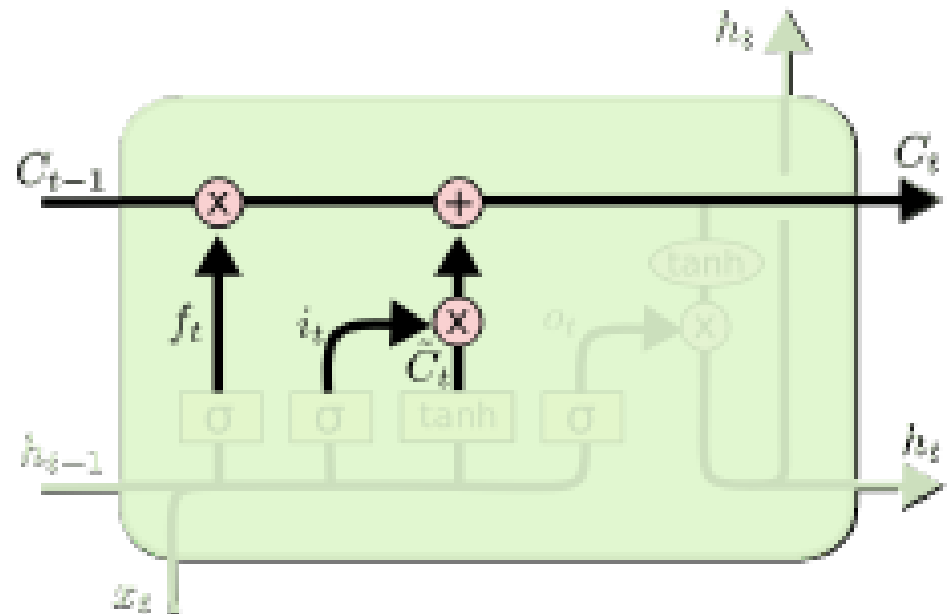$$i^{(t)} = \sigma(\mathbf{u}_i^\mathsf{T} \mathbf{x}^{(t)} + \mathbf{w}_i \mathbf{h}^{(t-1)})$$

$$\tilde{c}^{(t)} = \tanh(\mathbf{u}_{\tilde{c}}^\mathsf{T} \mathbf{x}^{(t)} + \mathbf{w}_{\tilde{c}} \mathbf{h}^{(t-1)})$$

# Step-by-Step LSTM Walk Through

▸ Third step: updates the old cell state $c^{(t-1)}$, into the new cell state $c^{(t)}$.

▸ We multiply the old state by $f^{(t)}$, forgetting the things we decided to forget earlier.

▸ Then we add it $i^{(t)} \otimes \tilde{c}^{(t)}$.

▸ This is the new candidate values, scaled by how much we decided to update each state value.

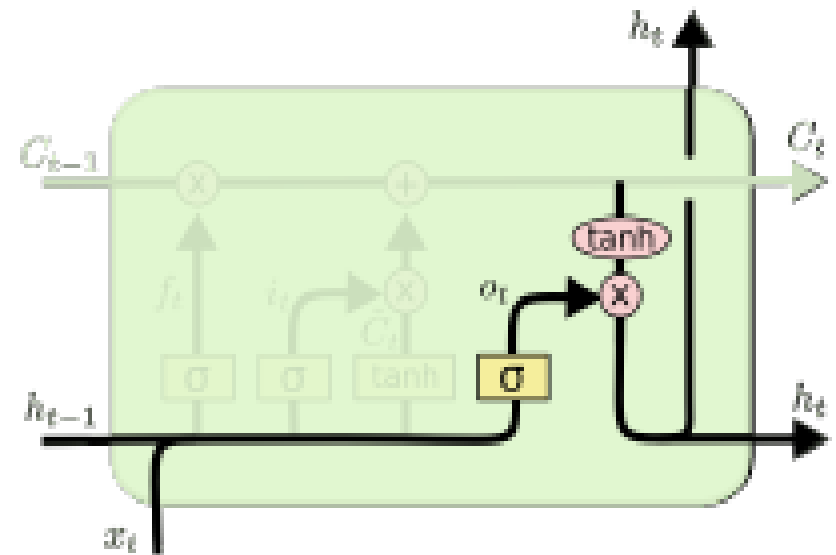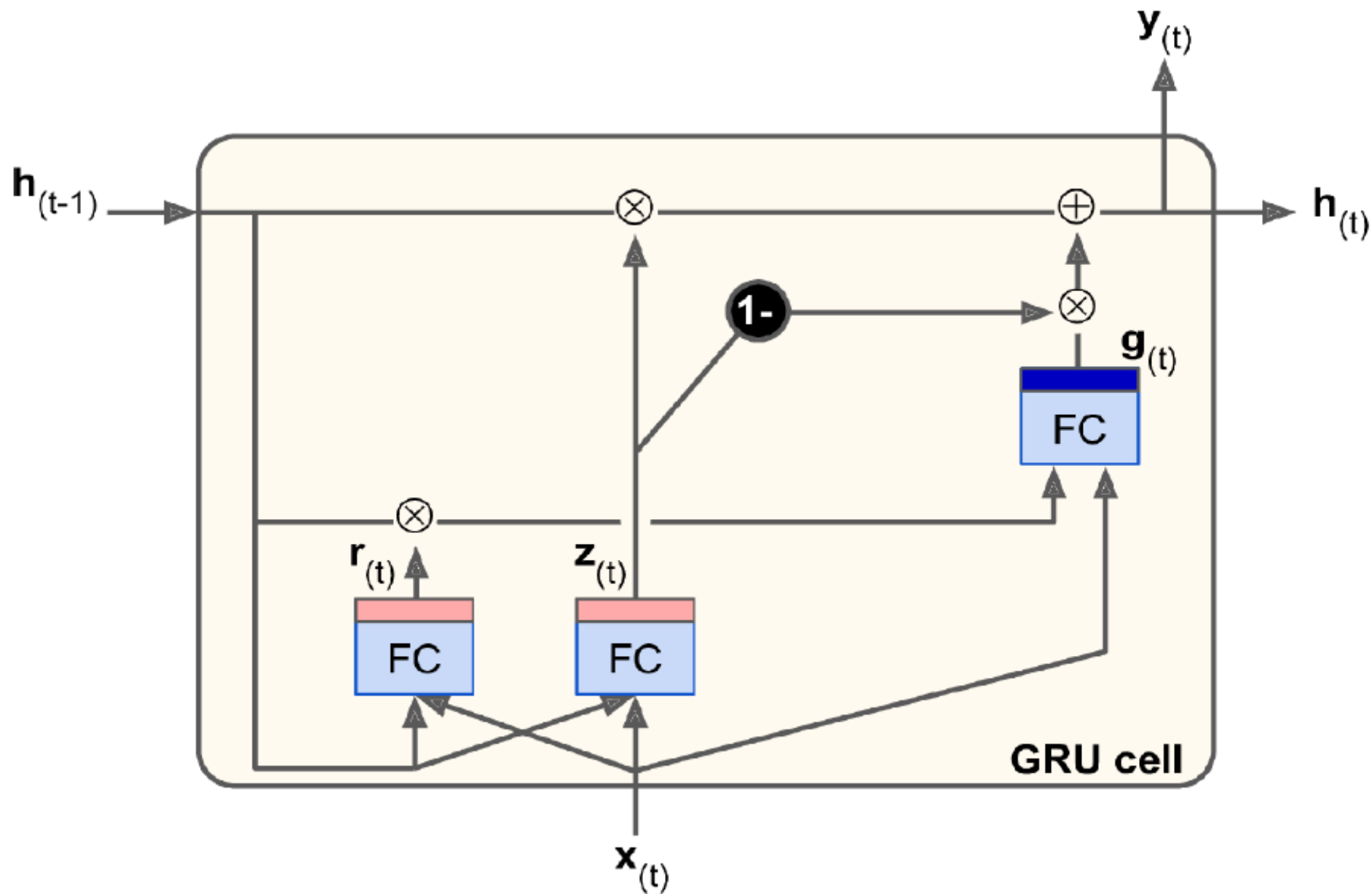$$c^{(t)} = f^{(t)} \otimes c^{(t-1)} + i^{(t)} \otimes \tilde{c}^{(t)}$$

# Step-by-Step LSTM Walk Through

▸ Fourth step: decides about the output.

▸ First, runs a sigmoid layer that decides what parts of the cell state we are going to output.

▸ Then, puts the cell state through tanh and multiplies it by the output of the sigmoid gate (output gate), so that it only outputs the parts it decided to.

$$o^{(t)} = \sigma(u_o^T x^{(t)} + w_o h^{(t-1)})$$

$$\hat{y}^{(t)} = h^{(t)} = o^{(t)} \otimes \tanh(c^{(t)})$$

# Gated Recurrent Unit (GRU)

# Gated Recurrent Unit (GRU)

$$\mathbf{z}_{(t)} = \sigma\left(\mathbf{W}_{xz}{}^{\top}\mathbf{x}_{(t)} + \mathbf{W}_{hz}{}^{\top}\mathbf{h}_{(t-1)} + \mathbf{b}_{z}\right)$$

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

$$\mathbf{g}_{(t)} = \tanh\left(\mathbf{W}_{xg}{}^{\top}\mathbf{x}_{(t)} + \mathbf{W}_{hg}{}^{\top}\left(\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}\right) + \mathbf{b}_{g}\right)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + \left(1 - \mathbf{z}_{(t)}\right) \otimes \mathbf{g}_{(t)}$$

# WaveNet- Use 1D convolutional layers



Conv1D dilation 8

Conv1D dilation 4

Conv1D dilation 2

Conv1D dilation 1

Input