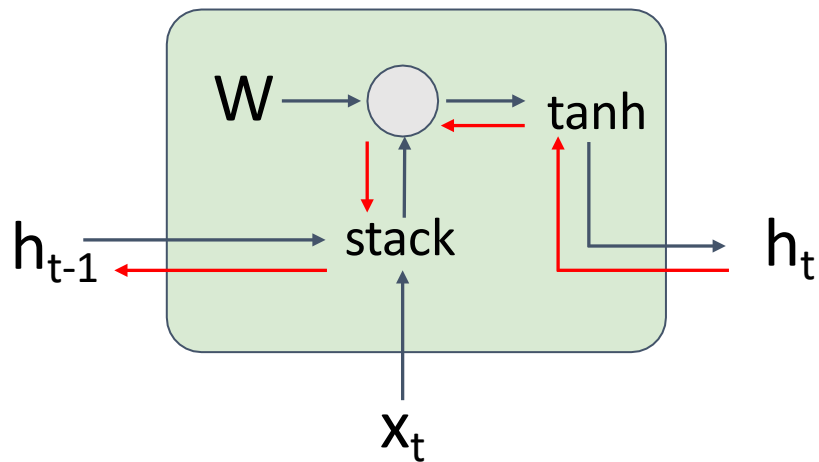# Problems learning RNNs

- In theory, should be able to propagate information over arbitrarily long contexts.

- In practice, RNNs suffer from **vanishing gradients** that decay to 0, or **exploding gradients** that increase towards infinity.
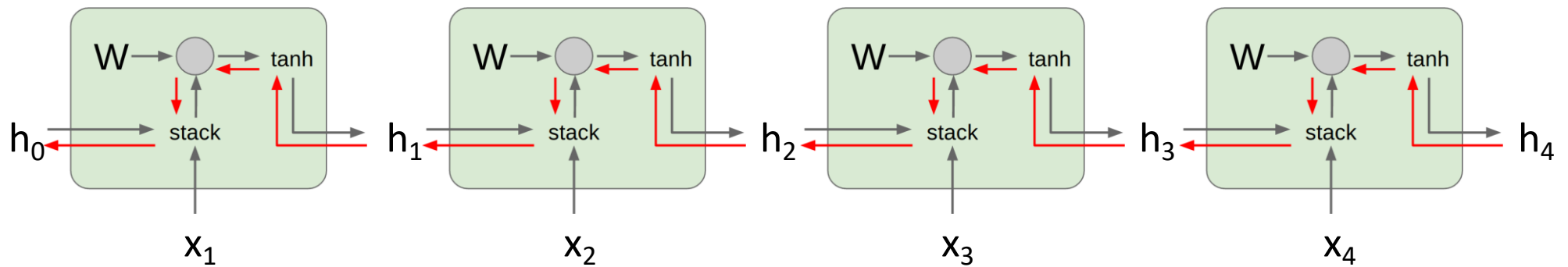
Figure: Goodfellow, Bengio and Courville. Deep Learning. MIT Press, 2016.

# Vanilla RNN Gradient Flow

Backpropagation from $h_t$ to $h_{t-1}$ multiplies by W
(actually $W_{hh}^T$)

W $\longrightarrow$ ◯ $\longrightarrow$ tanh

$h_{t-1}$ $\longrightarrow$ stack $\longrightarrow$ $h_t$

$x_t$

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

$$= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

$$= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$
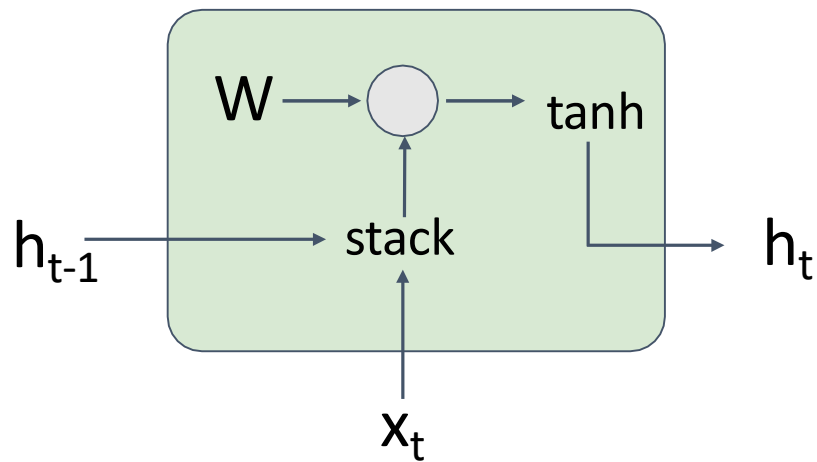
# Vanilla RNN Gradient Flow



Computing gradient of $h_0$
involves many  factors of W
(and repeated tanh)

# LSTM Networks

- Recurrent Neural Networks

  - ***Long Short-Term Memory (LSTM)*** networks are a variant of RNNs
  - LSTM mitigates the vanishing/exploding gradient problem
    - Solution: a Memory Cell, updated at each step in the sequence

  - Three gates control the flow of information to and from the Memory Cell
    - Input Gate: protects the current step from irrelevant inputs
    - Output Gate: prevents current step from passing irrelevant information to later steps
    - Forget Gate: limits information passed from one cell to the next

  - Most modern RNN models use either LSTM units or other more advanced types of recurrent units (e.g., GRU units)

# Vanilla RNN Gradient Flow



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

$$= \tanh\left( \begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

$$= \tanh\left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

# Long Short Term Memory (LSTM)

## Vanilla RNN

$$h_t = \tanh\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

## LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

Two vectors at each timestep:

Cell state

Hidden state

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM)

## Vanilla RNN

$$h_t = \tanh\left(W\begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

## LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

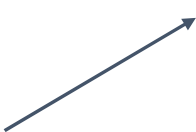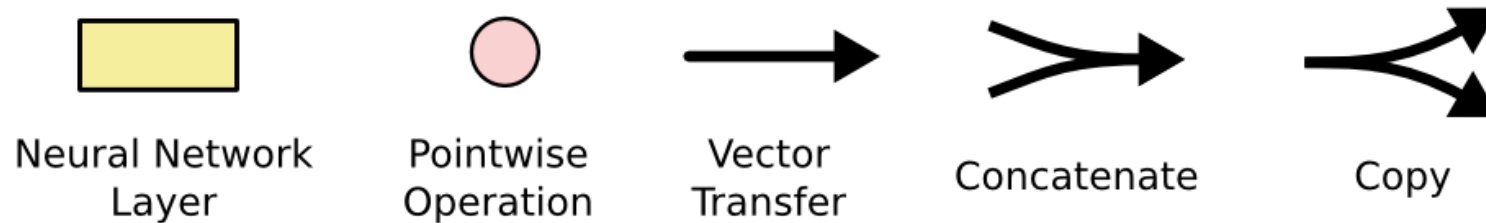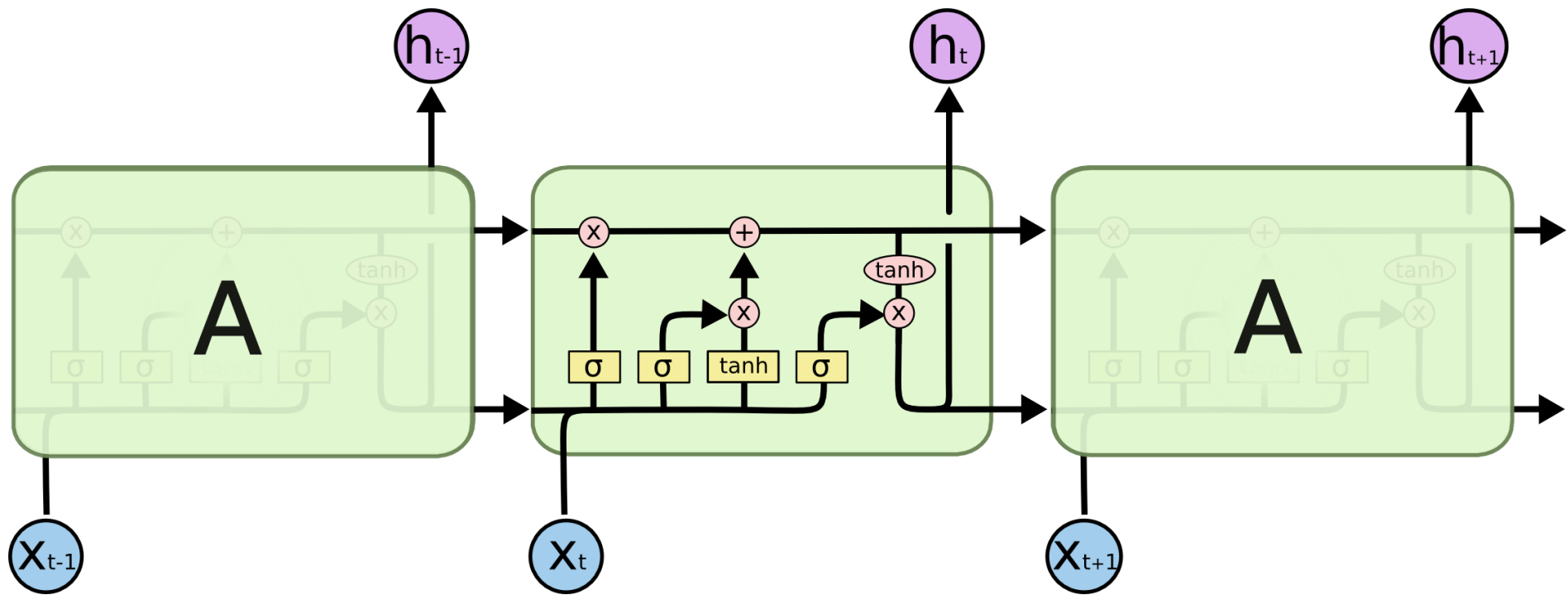Compute four **gates**
at each timestep
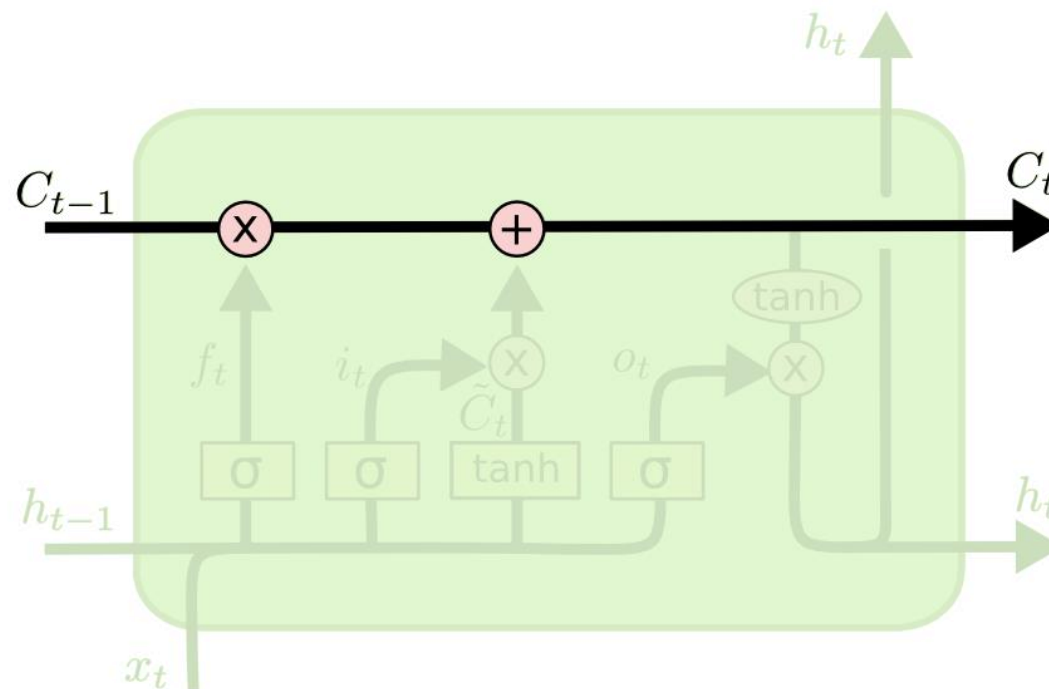
# Cell State Example

- Want to remember person & number of a subject noun so that it can be checked to agree with the person & number of verb when it is eventually encountered.

- Forget gate will remove existing information of a prior subject when a new one is encountered.

- Input gate "adds" in the information for the new subject.
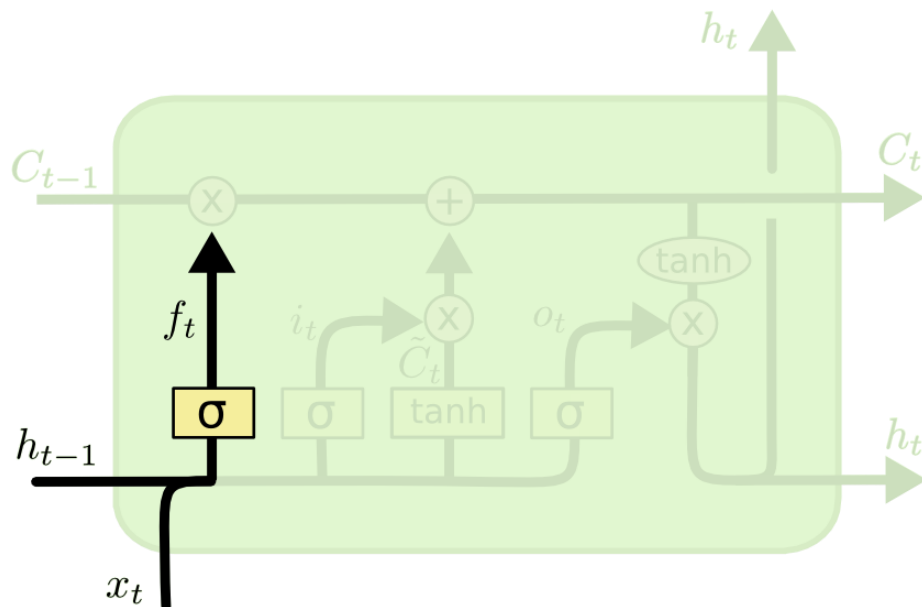
# LSTM Network Architecture

# Cell State

- Maintains a vector $C_t$ that is the same dimensionality as the hidden state, $h_t$
- Information can be added or deleted from this state vector via the forget and input gates.
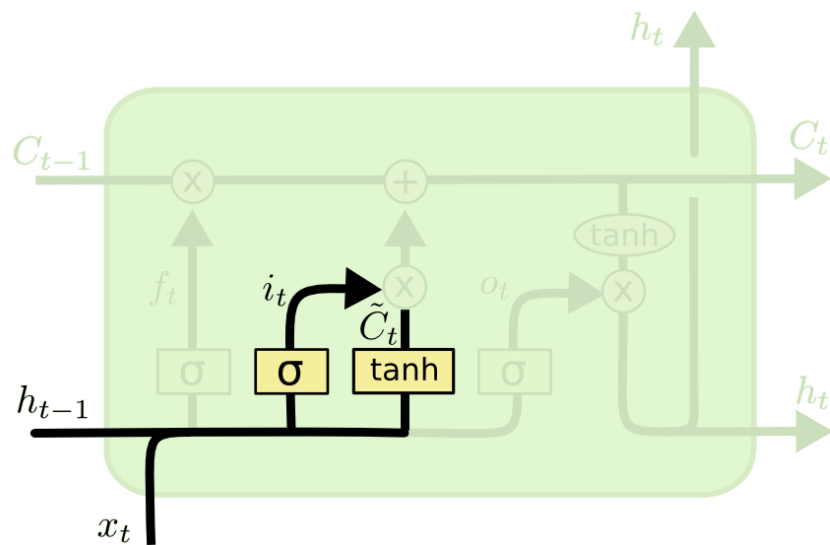
# Forget Gate

- Forget gate computes a 0-1 value using a logistic sigmoid output function from the input, $x_t$, and the current hidden state, $h_t$:

- Multiplicatively combined with cell state, "forgetting" information where the gate outputs something close to 0.

$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

11

# Input Gate

- First, determine which entries in the cell state to update by computing 0-1 sigmoid output.

- Then determine what amount to add/subtract from these entries by computing a tanh output (valued −1 to 1) function of the input and hidden state.
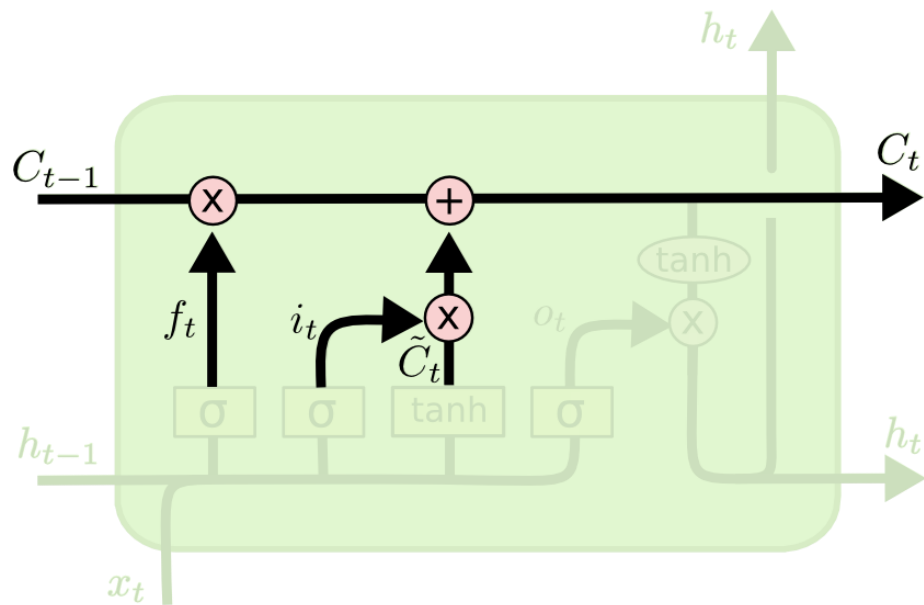
$$i_t = \sigma \left( W_i \cdot [h_{t-1}, x_t] \; + \; b_i \right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] \; + \; b_C)$$
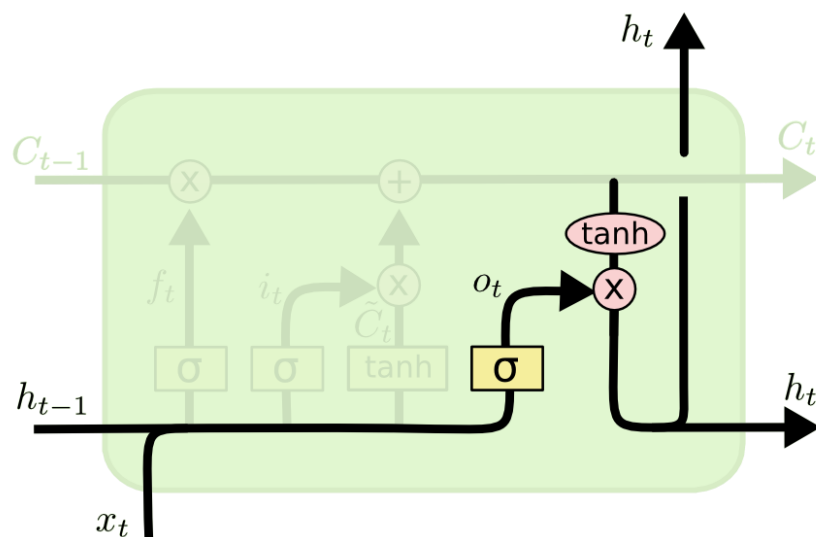
# Updating the Cell State

- Cell state is updated by using component-wise vector multiply to "forget" and vector addition to "input" new information.

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# Output Gate

- Hidden state is updated based on a "filtered" version of the cell state, scaled to −1 to 1 using tanh.

- Output gate computes a sigmoid function of the input and current hidden state to determine which elements of the cell state to "output".
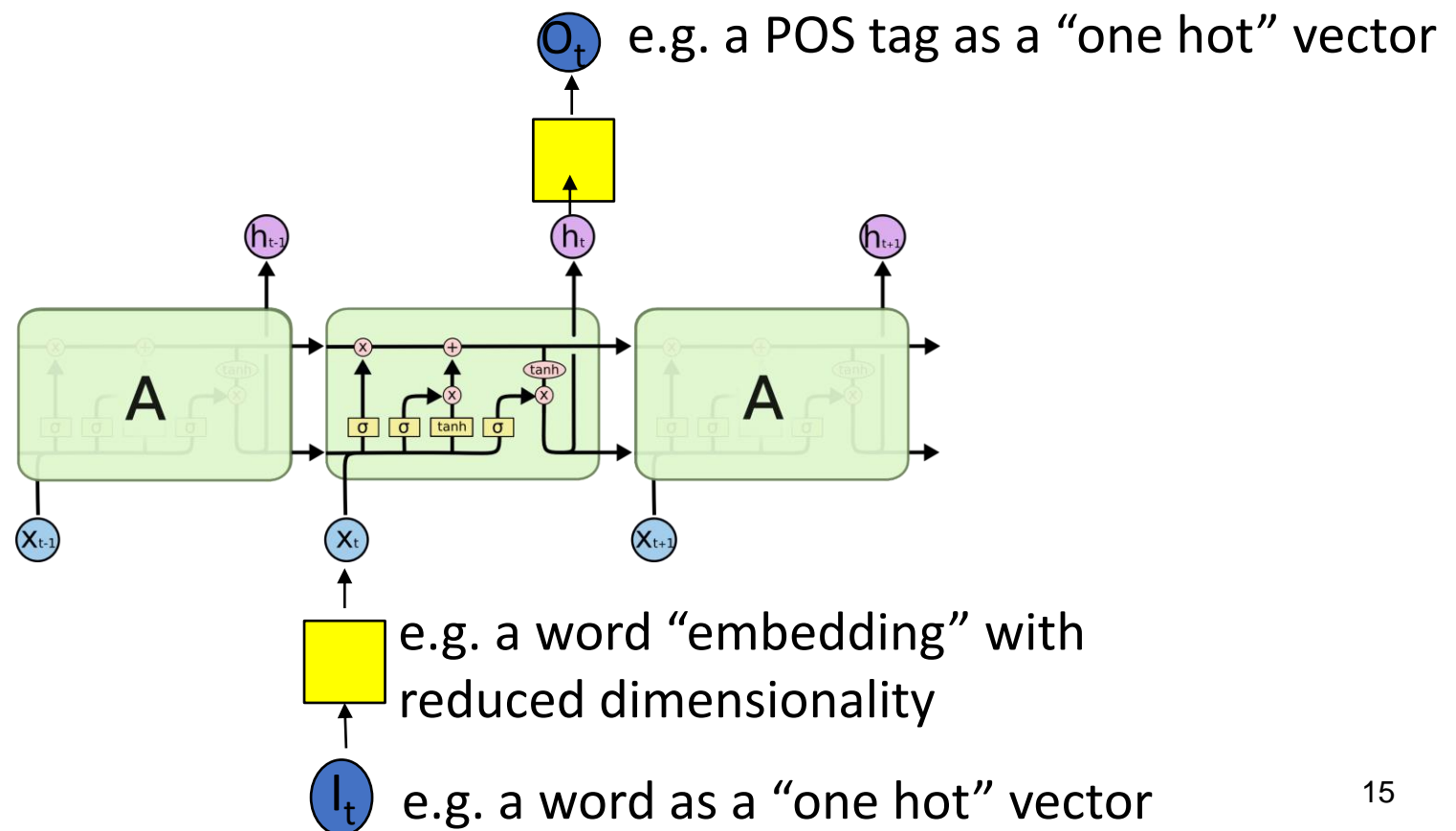
$$o_t = \sigma\left(W_o\left[h_{t-1}, x_t\right] + b_o\right)$$

$$h_t = o_t * \tanh\left(C_t\right)$$

14

# Overall Network Architecture

- Single or multilayer networks can compute LSTM inputs from problem inputs and problem outputs from LSTM outputs.

$O_t$ e.g. a POS tag as a "one hot" vector

e.g. a word "embedding" with reduced dimensionality

$I_t$ e.g. a word as a "one hot" vector

# LSTM Training

- Trainable with backprop derivatives such as:
  - Stochastic gradient descent (randomize order of examples in each epoch) with momentum (bias weight changes to continue in same direction as last update).
  - ADAM optimizer (Kingma & Ma, 2015)
- Each cell has many parameters ($W_f$, $W_i$, $W_C$, $W_o$)
  - Generally requires lots of training data.
  - Requires lots of compute time that exploits GPU clusters.