# Introduction to Big Data Science

## 12-3 Period

## Understanding Long-Short Term Memory (LSTM) Networks

# Contents

- **Sequence Data and RNN**
- **Problem of Long Term Dependency**
- **LSTM Networks**
- **Core Idea behind LSTM**
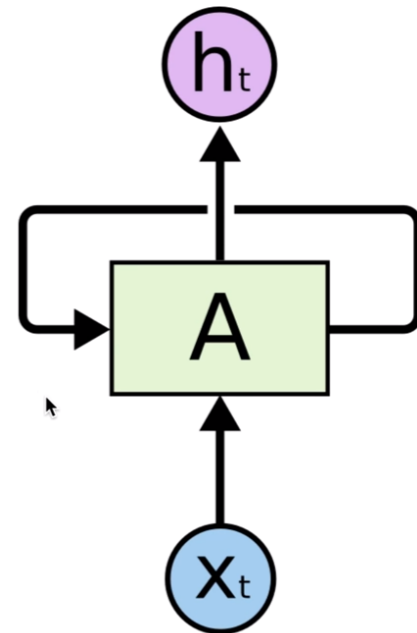- **Working Steps of LSTM**
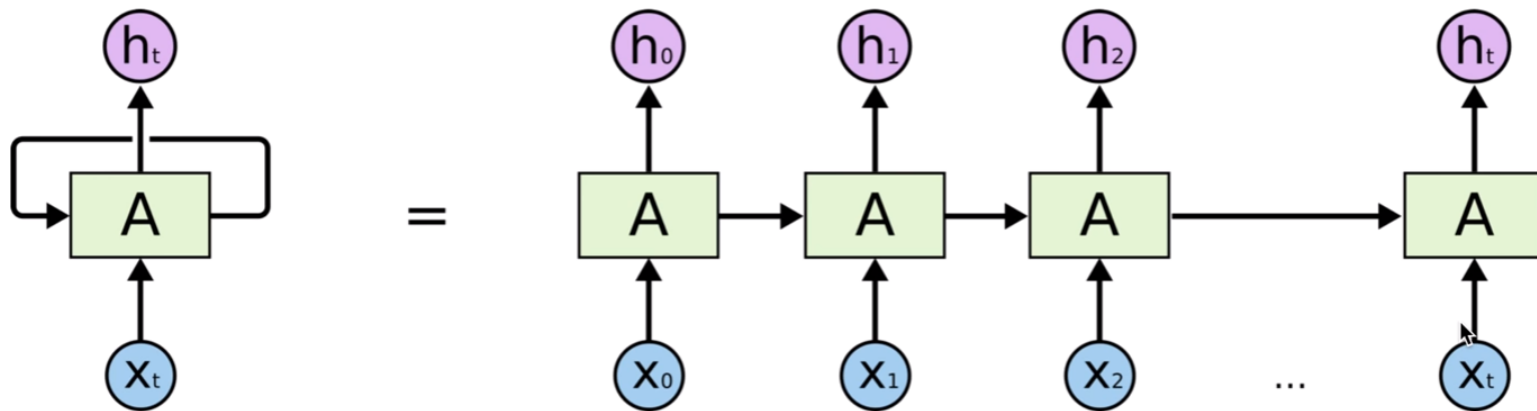- **Variants on LSTM**

# Sequence Data

## Sequence data

- We don't understand one word only

- We understand based on the previous words + this word. (time series)
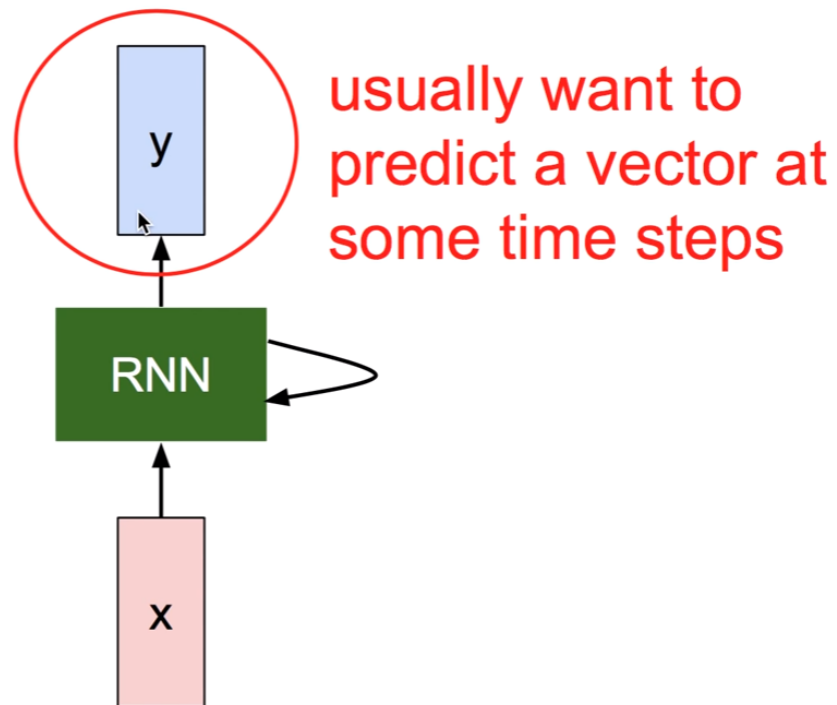
- NN/CNN cannot do this

# Sequence Data

- We don't understand one word only

- We understand based on the previous words + this word. (time series)

- NN/CNN cannot do this



http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Recurrent Neural Network



usually want to predict a vector at some time steps

# Recurrent Neural Network

We can process a sequence of vectors **x** by
applying a recurrence formula at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$
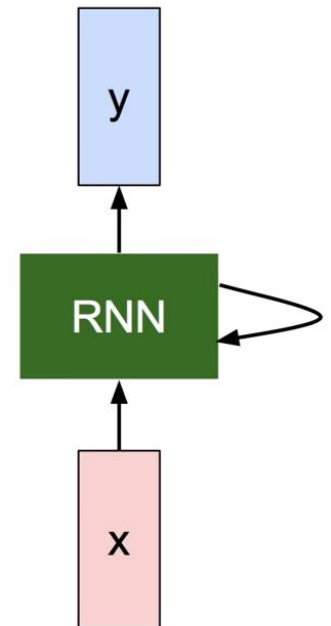
new state

old state  input vector at
some time step

some function
with parameters W

y

RNN

x

# (Vanilla) Recurrent Neural Network

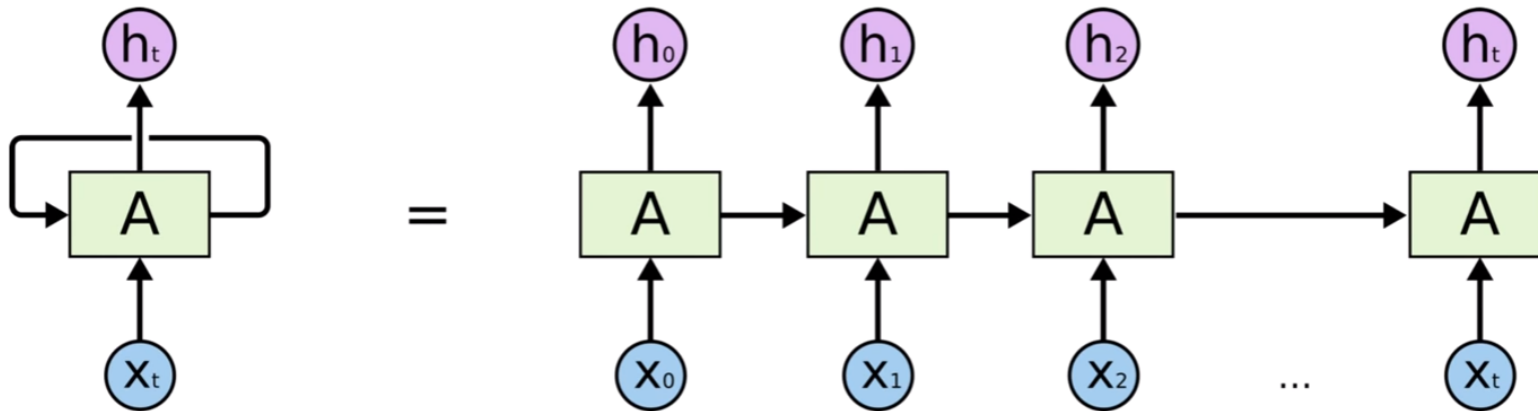The state consists of a single *"hidden"* vector **h**:

$$h_t = f_W(h_{t-1}, x_t)$$

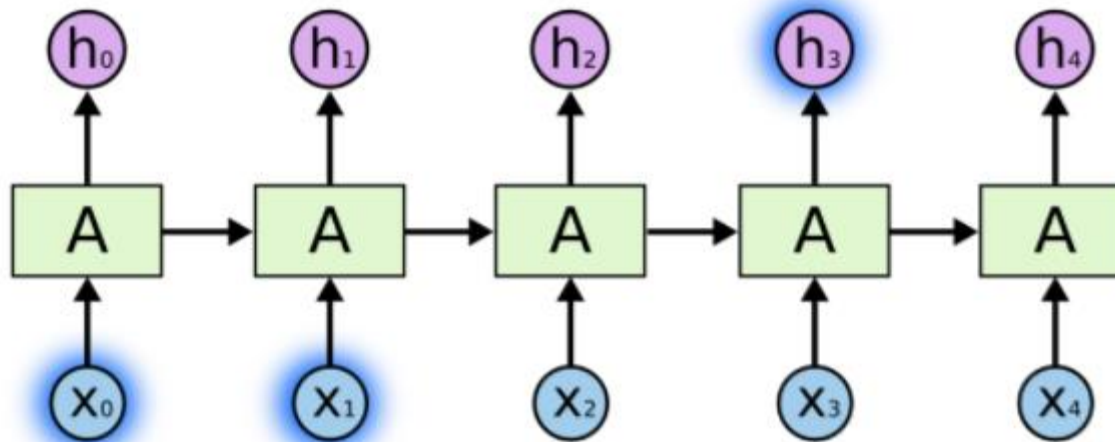$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$y_t = W_{hy} h_t$$

Notice: the same function and the same set of parameters are used at every time step.
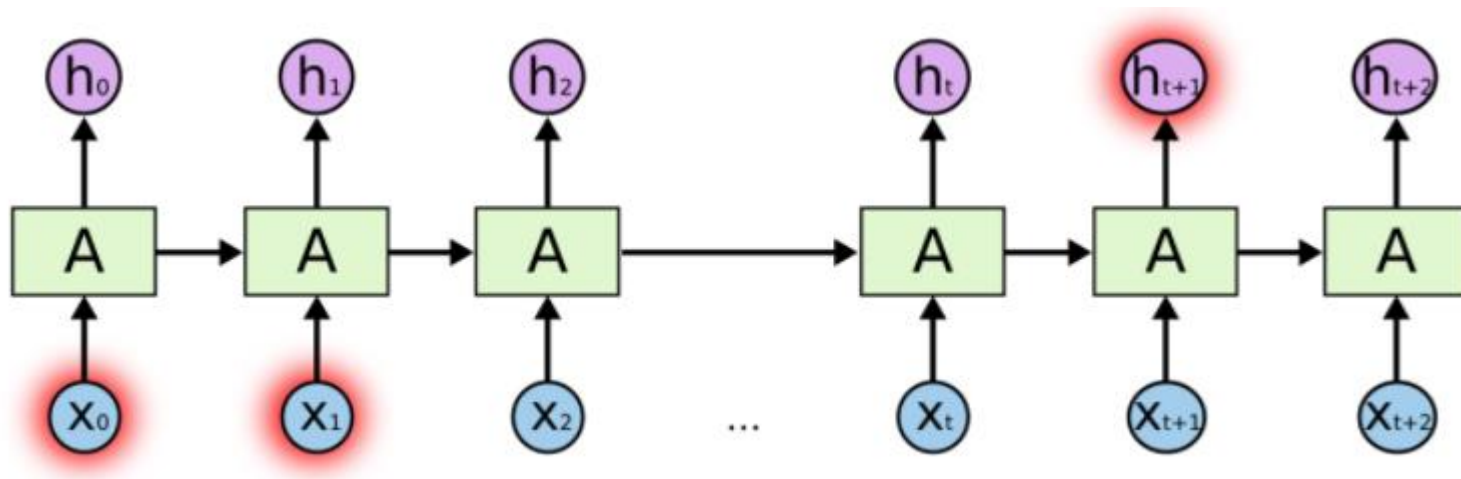
http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Problem of Long Term Dependencies

◆ If we are trying to predict the last word in "the clouds are in the *sky*," we don't need any further context – it's pretty obvious the next word is going to be sky.

◆ 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.
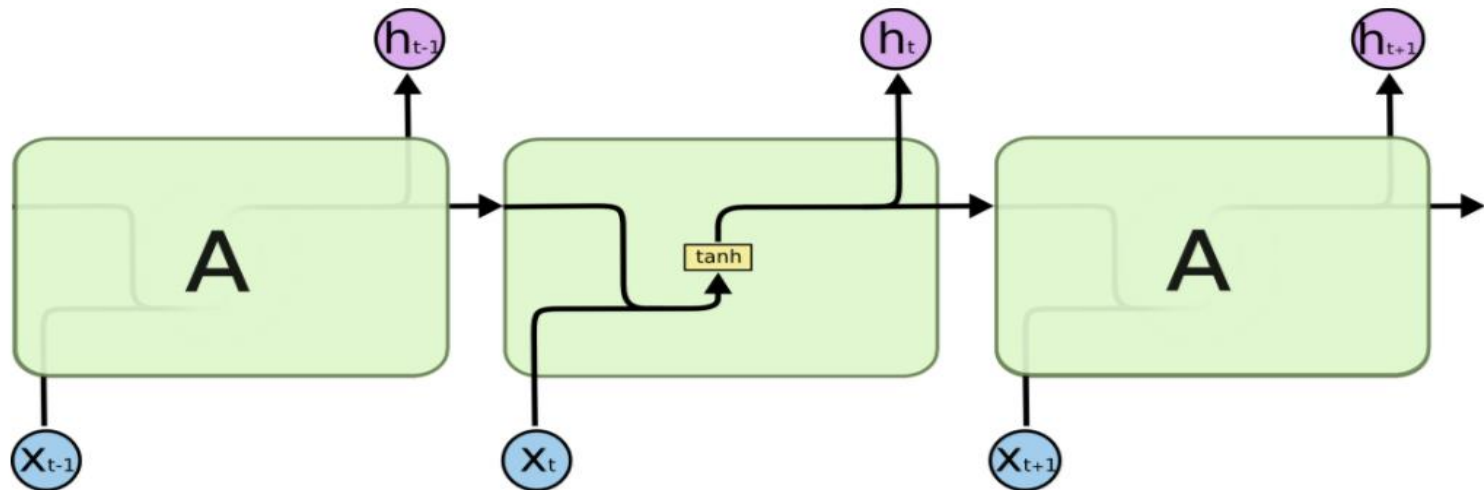
# Problem of Long Term Dependencies

◆ An example: To predict the last word in the text "I grew up in France… I speak fluent *French*." Recent information suggests that the next word is probably the name of a language, but if we want to narrow down which language, **we need the context of France, from further back**.

◆ It's entirely possible for the gap between the relevant information and the point where it is needed to become very large.

◆ Unfortunately, as that gap grows, RNNs become unable to learn to connect the information!
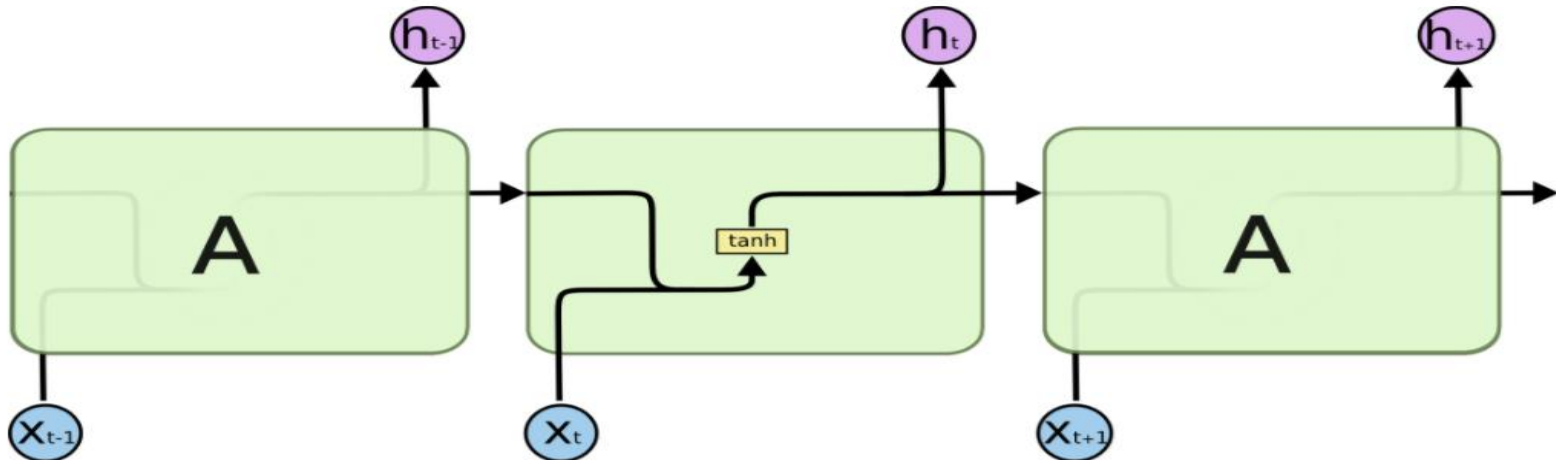
# LSTM Network

◆ Long Short Term Memory networks – usually just called "LSTMs" – are a special kind of RNN, capable of learning long-term dependencies. They were introduced by Hochreiter & Schmidhuber (1997), and were refined and popularized by many people in following work.1 They work tremendously well on a large variety of problems, and are now widely used.



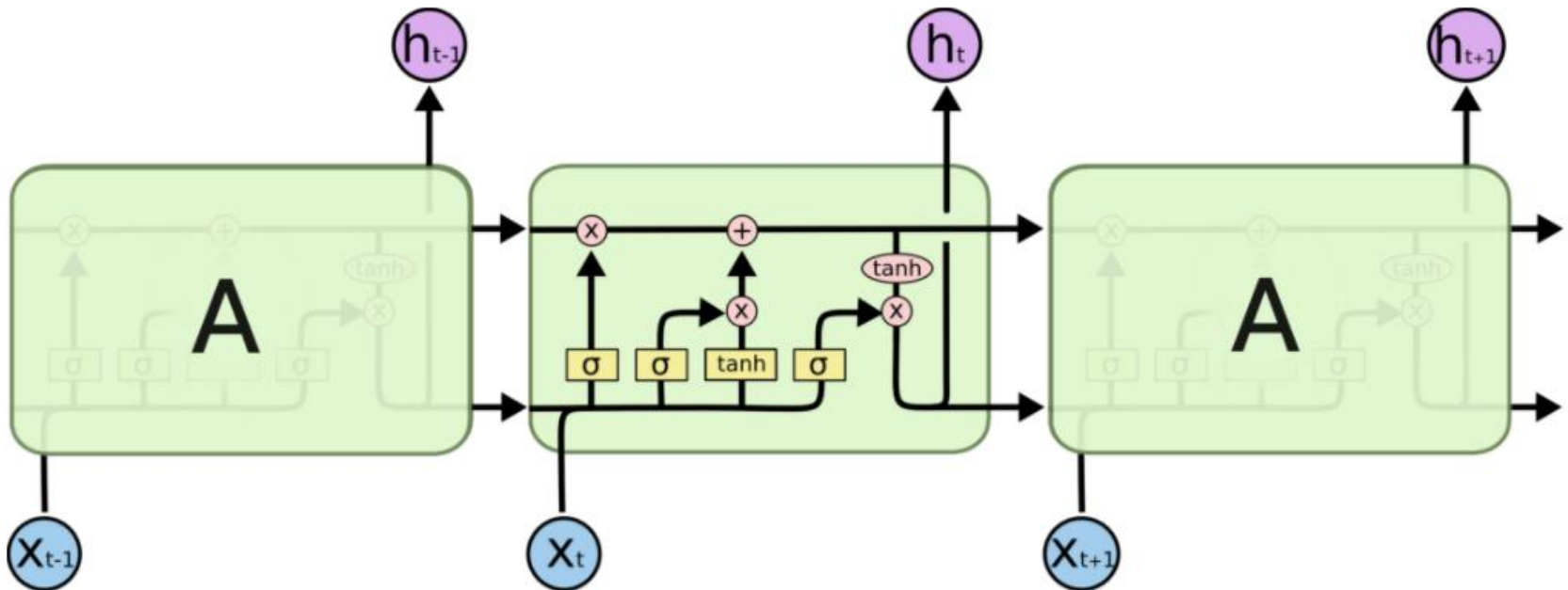The repeating module in a standard RNN contains a single layer.

# LSTM Network

◆ LSTMs are explicitly designed to avoid the long-term dependency problem. Remembering information for long periods of time is practically their default behavior, not something they struggle to learn!

◆ All recurrent neural networks have the form of a chain of repeating modules of neural network. In standard RNNs, this repeating module will have a very simple structure, such as a single tanh layer.
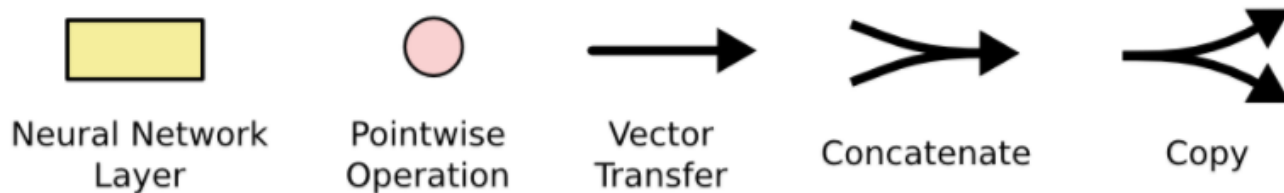


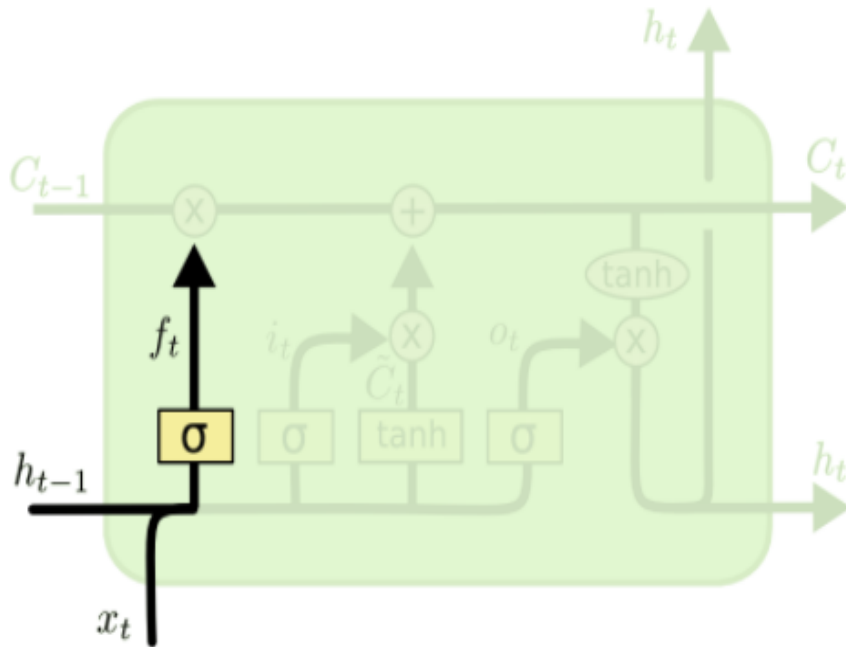The repeating module in a standard RNN contains a single layer.

# LSTM Network



The repeating module in an LSTM contains four interacting layers.
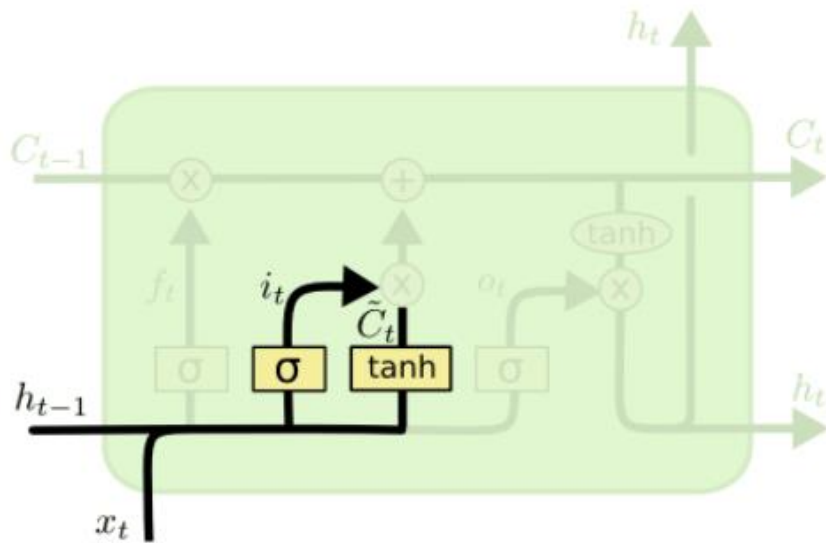
Neural Network Layer · Pointwise Operation · Vector Transfer · Concatenate · Copy

# Working Steps of LSTM



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

# Working Steps of LSTM



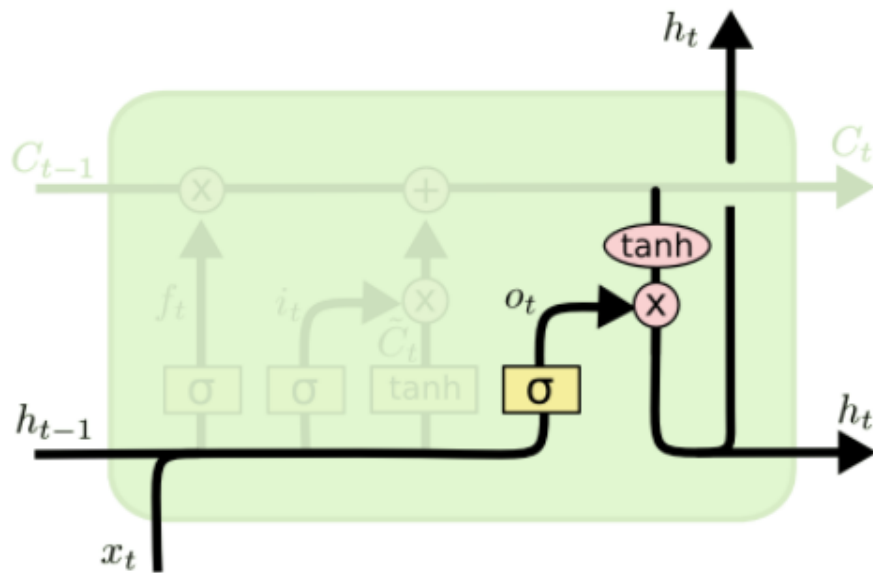$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] \; + \; b_i\right)$$

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

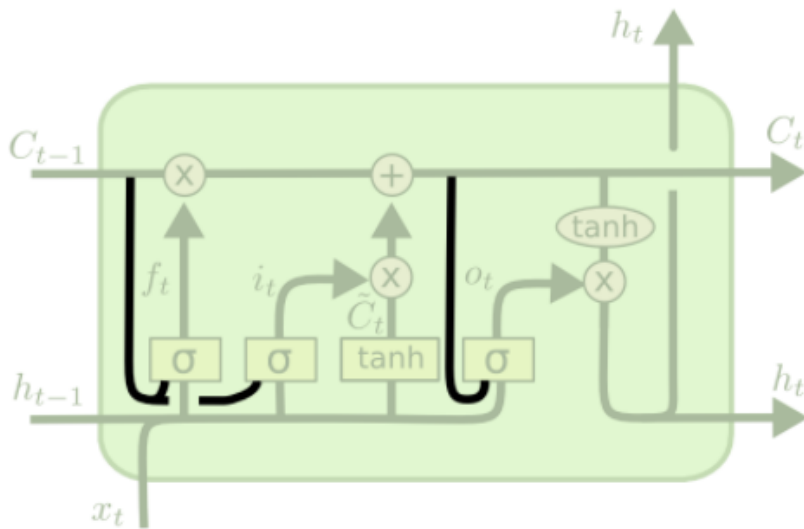# Working Steps of LSTM

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

# Working Steps of LSTM



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

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

# Variants on LSTM



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

$$i_t = \sigma\left(W_i \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_t] \; + \; b_i\right)$$

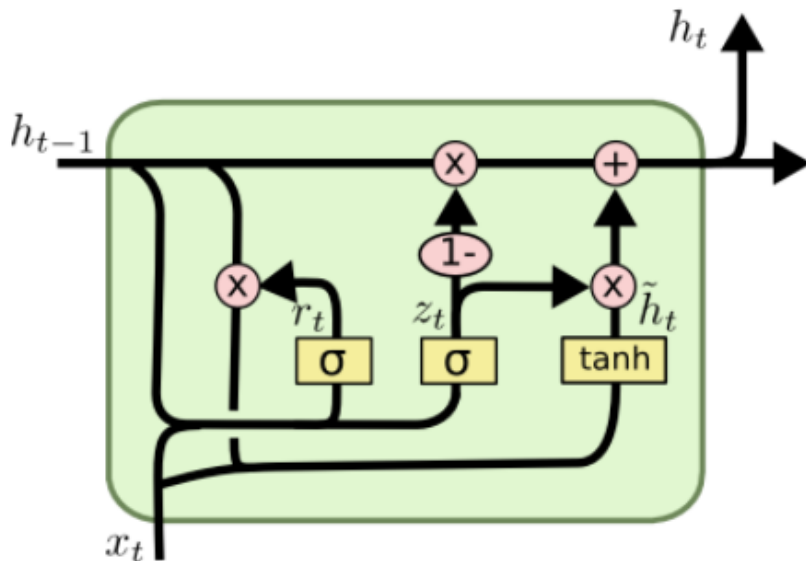$$o_t = \sigma\left(W_o \cdot [\boldsymbol{C_t}, h_{t-1}, x_t] \; + \; b_o\right)$$

# Variants on LSTM



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

# Variants on LSTM

◆A slightly more dramatic variation on the LSTM is the Gated Recurrent Unit, or GRU, introduced by Cho, et al. (2014).
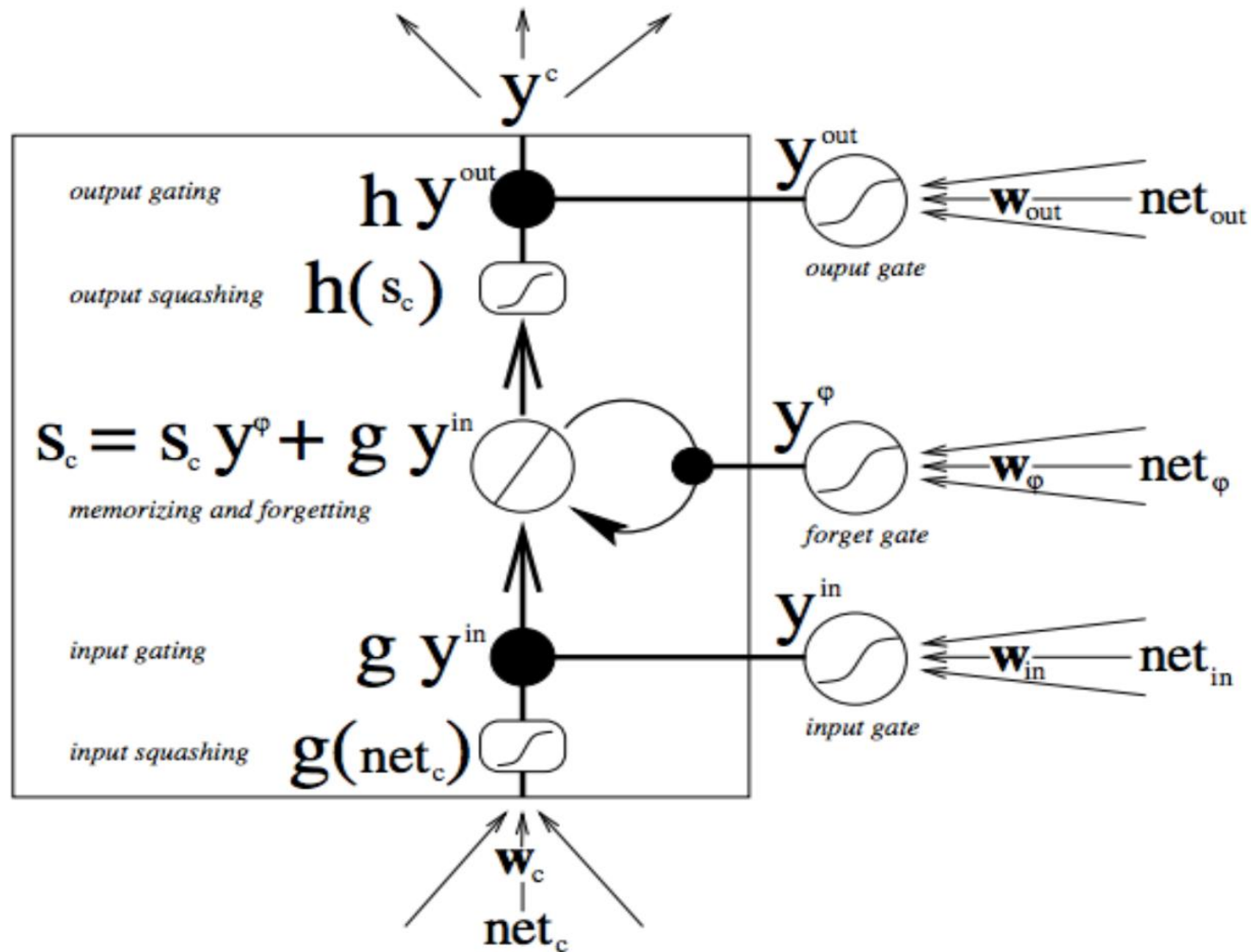


$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Additional Explanation on LSTM
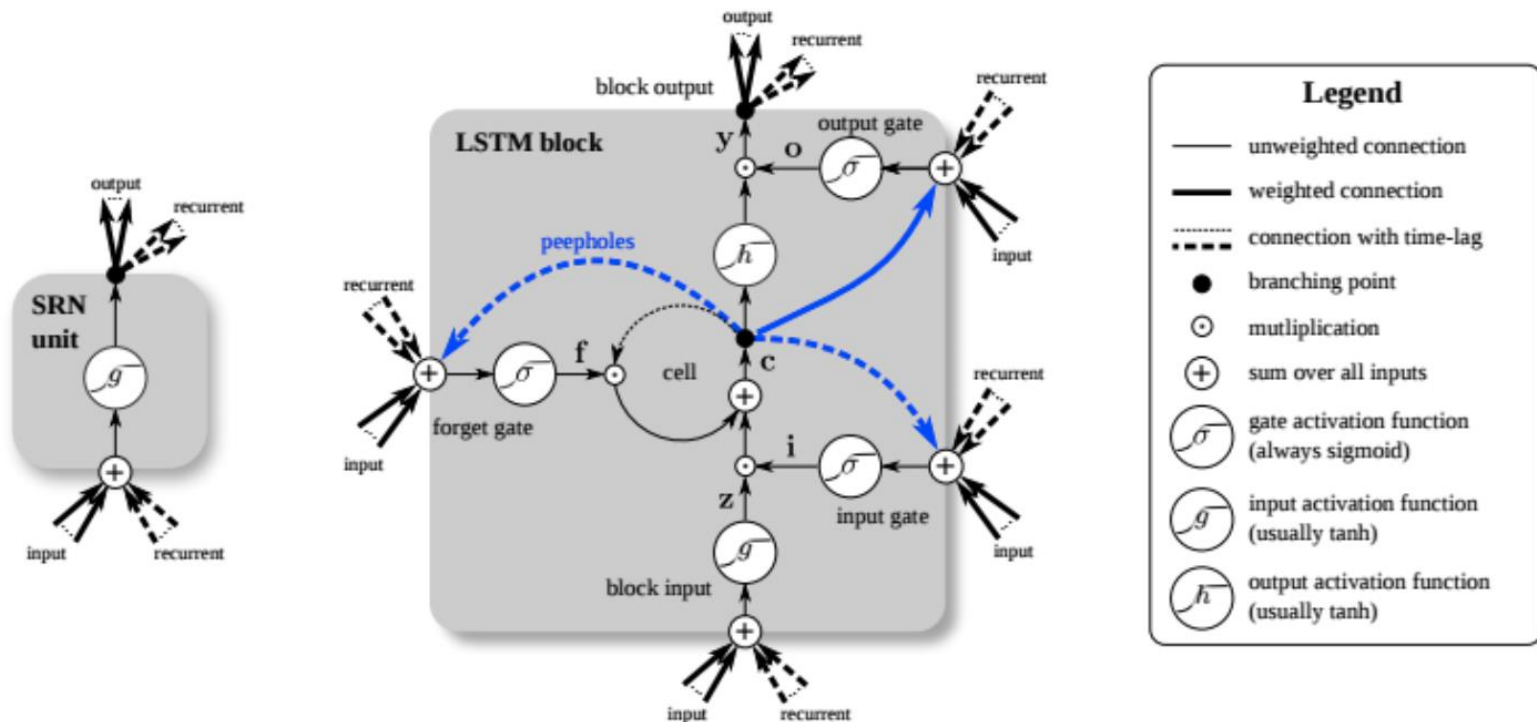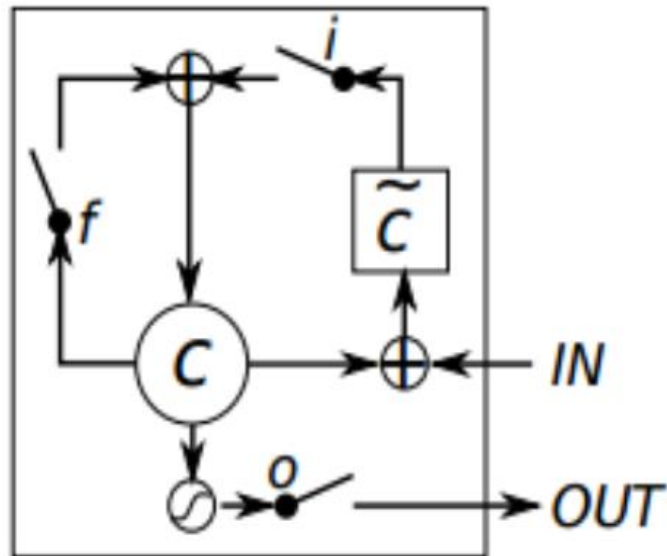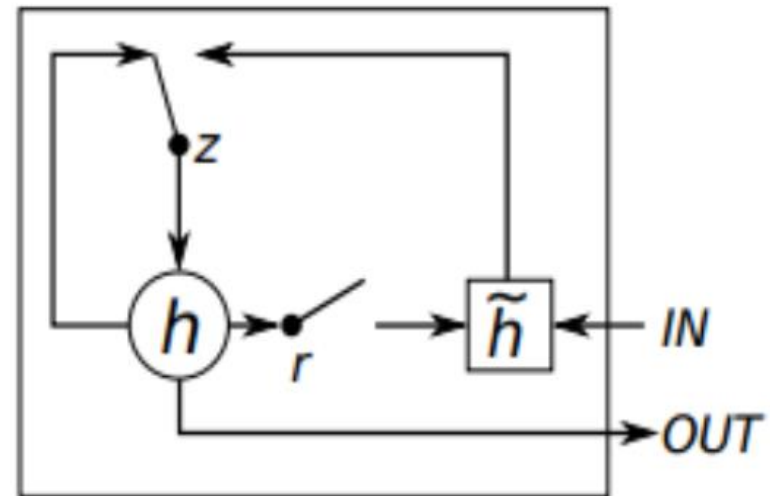
# Additional Explanation on LSTM



Figure 1. Detailed schematic of the Simple Recurrent Network (SRN) unit (left) and a Long Short-Term Memory block (right) as used in the hidden layers of a recurrent neural network.

# Additional Explanation on LSTM



(a) Long Short-Term Memory

(b) Gated Recurrent Unit

Figure 1: Illustration of (a) LSTM and (b) gated recurrent units. (a) $i$, $f$ and $o$ are the input, forget and output gates, respectively. $c$ and $\tilde{c}$ denote the memory cell and the new memory cell content. (b) $r$ and $z$ are the reset and update gates, and $h$ and $\bar{h}$ are the activation and the candidate activation.