# Chapter-3-part-II
# Recurrent Neural Networks

**https://www.geeksforgeeks.org/ml-back-propagation-through-time/**

<start> Giraffes standing <end>

Pretrained CNN using ImageNet dataset

CNN

Feature vector at fc layer (1×1×2048)

Linear

Softmax Softmax Softmax Softmax

LSTM LSTM LSTM — ... → LSTM

$W_{emb}$ $W_{emb}$ $W_{emb}$

<start> Giraffes other

$Y_T$ $y_{T-n}$ $y_t$ $y_{T-1}$

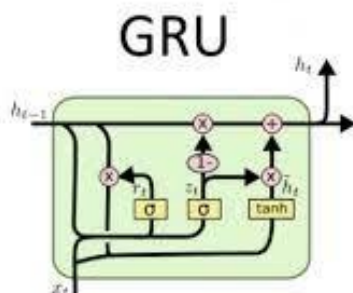$W_{hy}$ $W_{hy}$ $W_{hy}$

RNN = $h_{T-n}$ $W_{hh}$ $h_t$ $W_{hh}$ $h_{T-1}$

$W_{xh}$ $W_{xh}$ $W_{xh}$

$X_T$ $x_{T-n}$ $x_t$ $x_{T-1}$

1 Standard RNN architecture and an unfolded structure with T time

one to one   one to many   many to one   many to many   many to many

DECODER

መማር አዲስ መረዳትን የማግኘት ሂደት ነው <EOS>

$(h_0,c_0)$ $(h_1,c_1)$ $(h_2,c_2)$ $(h_3,c_3)$ $(h_4,c_4)$ $(h_5,c_5)$ $(h_6,c_6)$ $(h_7,c_7)$

<SOS> መማር አዲስ መረዳትን የማግኘት ሂደት ነው

ENCODER

$(h_0,c_0)$ $(h_1,c_1)$ $(h_2,c_2)$ $(h_3,c_3)$ $(h_4,c_4)$ $(h_5,c_5)$ $(h_6,c_6)$ $(h_7,c_7)$ $(h_8,c_8)$

learning is the process of gaining new understanding

LSTM GRU

እንደጻፈልን፣ ከሠላምታ ጋር

output text

CTC-decoder

LSTM

LSTM

BLSTM

$p_1$ $p_2$ $p_3$ ... $p_t$

Input image   Convolutional layers   Soft-max outputs

Transcription layer

Visual Representation

CNN

Combine

Predict answer

What is the color of the bird?

Textual Representation

White

Word /Sentence embedding+LSTM

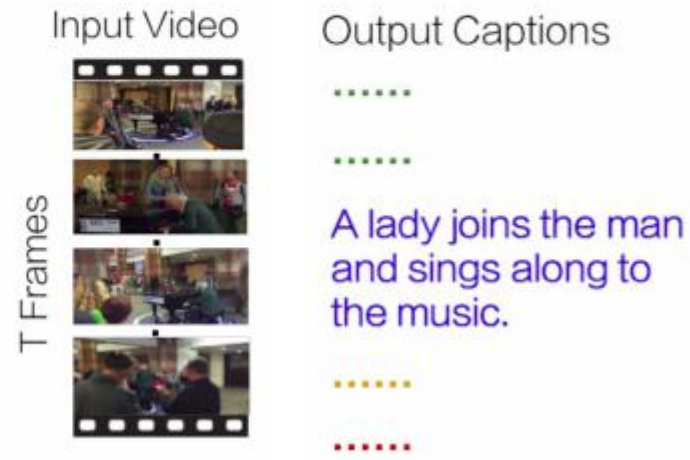# What is Recurrent Neural Networks?

## Why existing convnets are insufficient?

### Variable sequence length inputs and outputs!

**Example task**: video captioning

**Input** video can have variable number of frames

**Output** captions can be variable length.



Input Video     Output Captions

T Frames

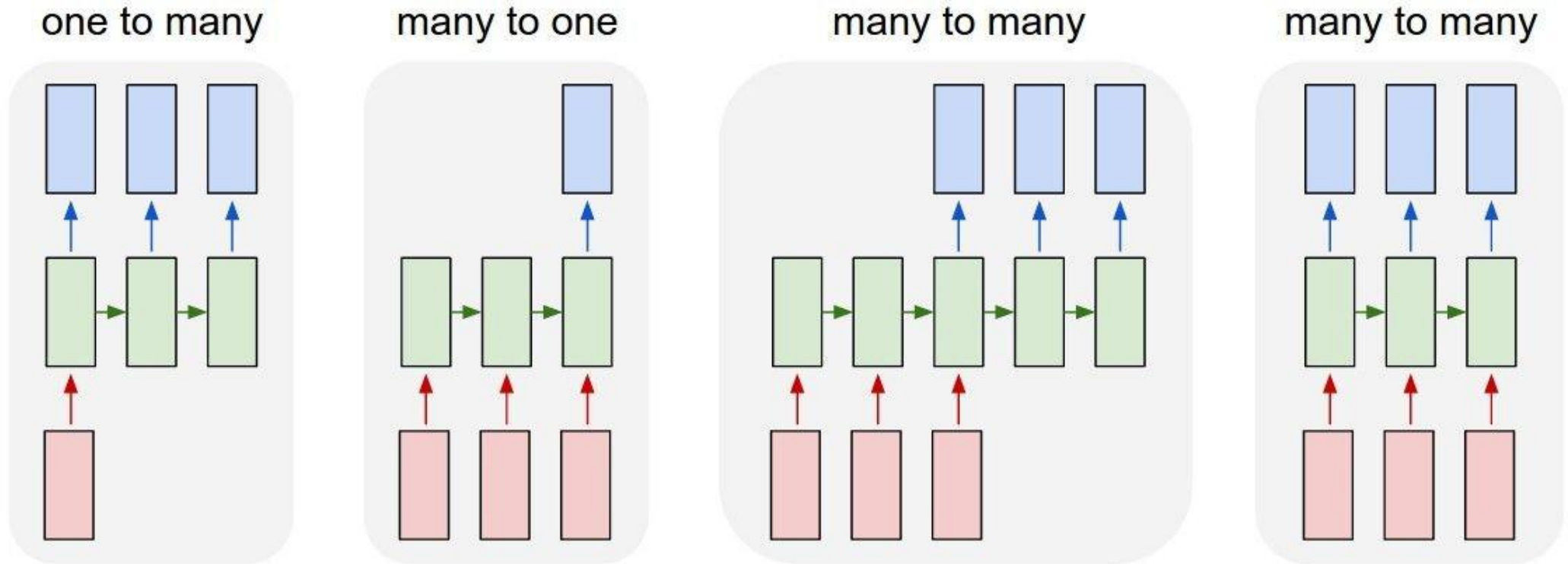A lady joins the man and sings along to the music.
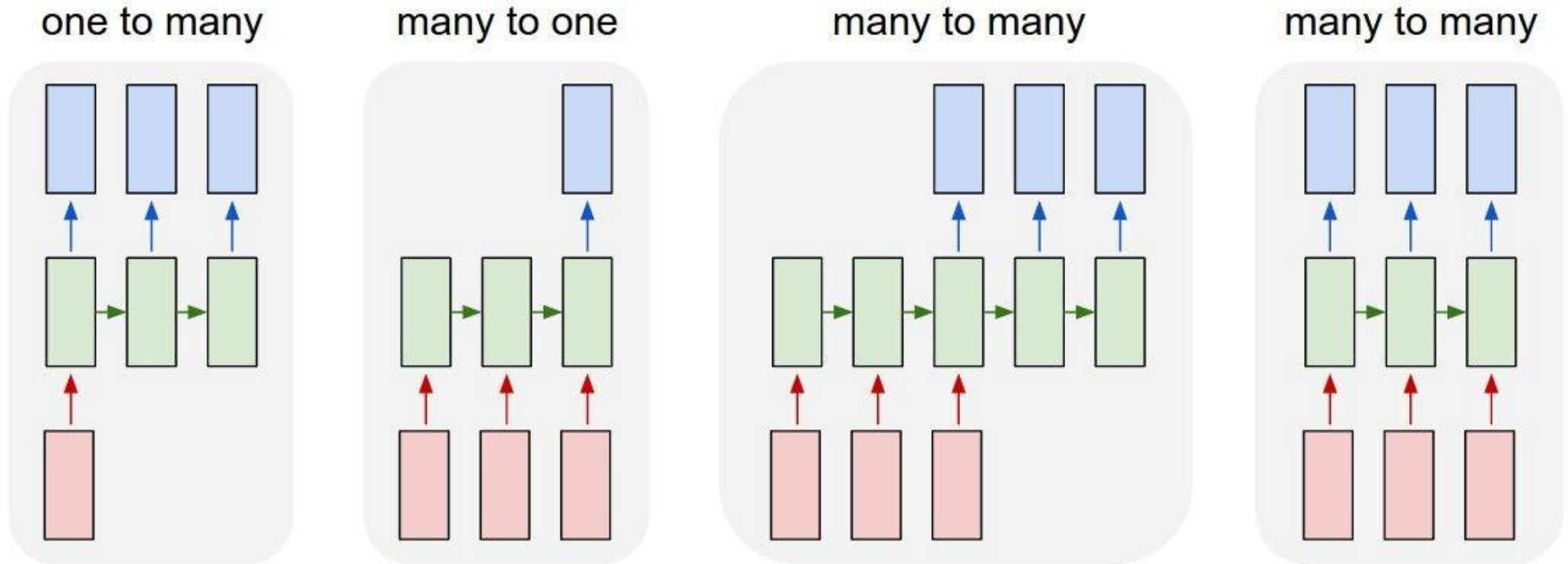
# What is Recurrent Neural Networks?

## Motivation

- In artificial intelligence applications we often deal with SEQUENCES (e.g. speech, video, machine translation…)
- We need models that take information of SEQUENCE


- HOW?
- By recurrency
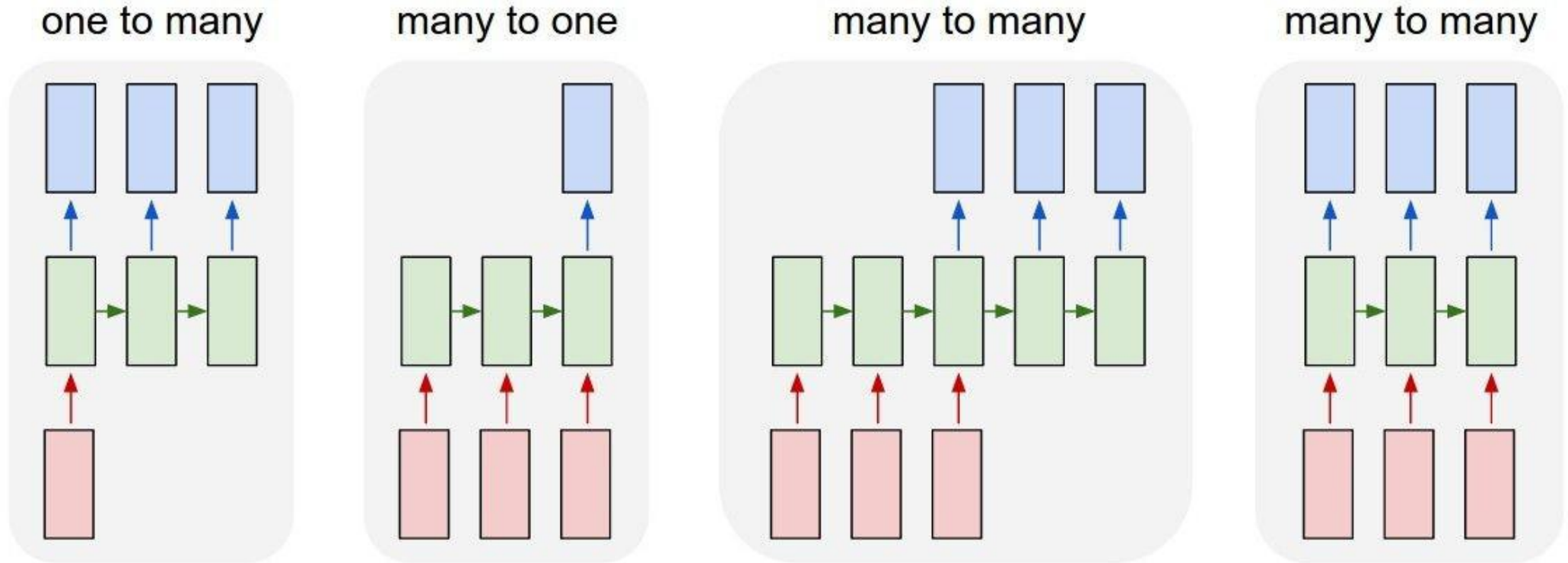
# Recurrent Neural Networks: Process Sequences

one to many

many to one

many to many

many to many

e.g. **Image Captioning**
image -> sequence of words

# Recurrent Neural Networks: Process Sequences



one to many    many to one    many to many    many to many
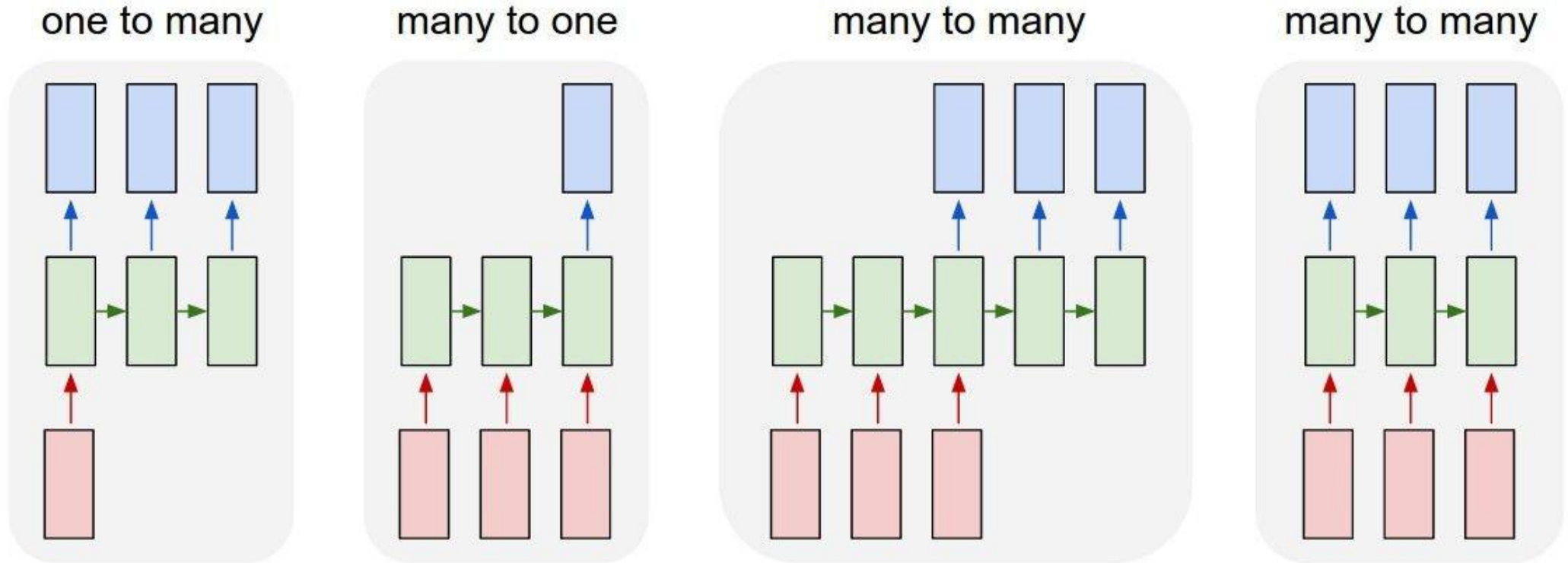
e.g. **Sentiment Classification**
sequence of words -> sentiment

# Recurrent Neural Networks: Process Sequences



one to many

many to one

many to many

many to many

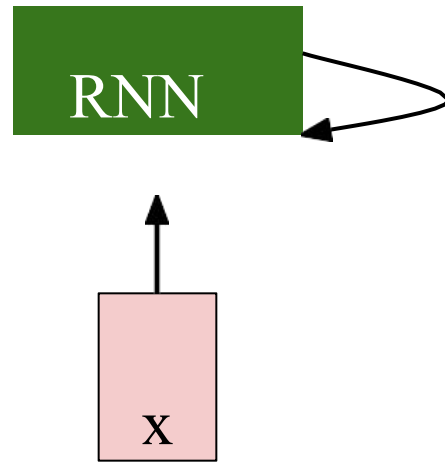e.g. **Machine Translation**
seq of words -> seq of words

# Recurrent Neural Networks: Process Sequences



e.g. **Video classification on frame level**

# Recurrent Neural Network

# 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

some function
with parameters W

old state input vector at
some time step

# Recurrent Neural Network...

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

$$h_t = f_W(h_{t-1}, x_t)$$
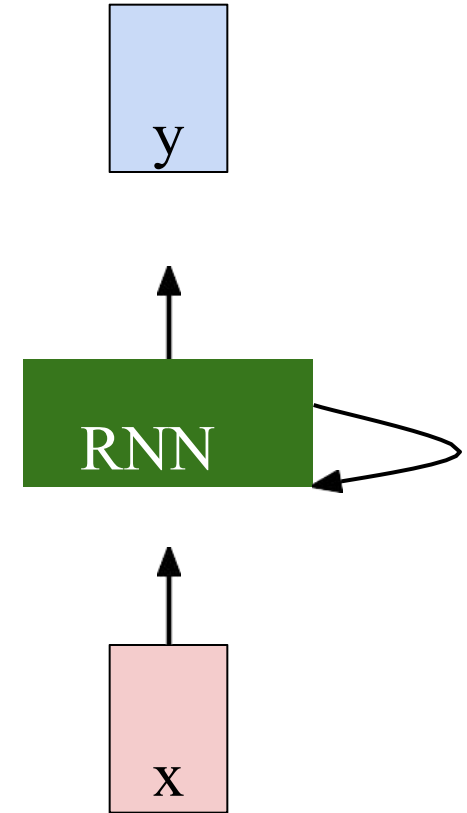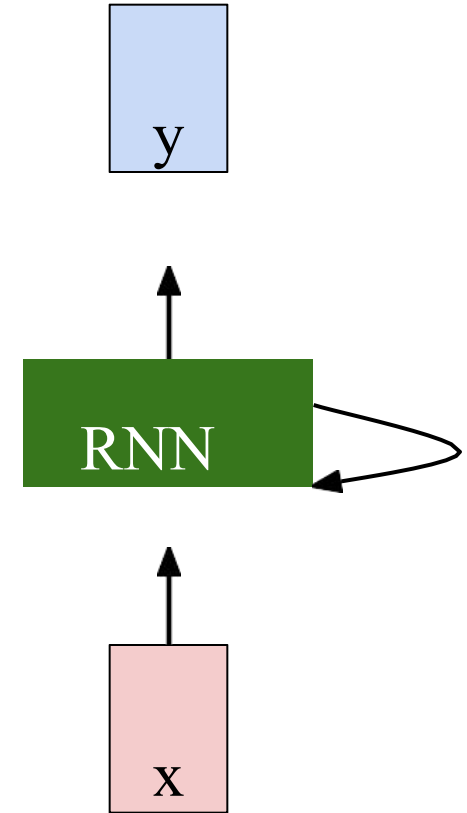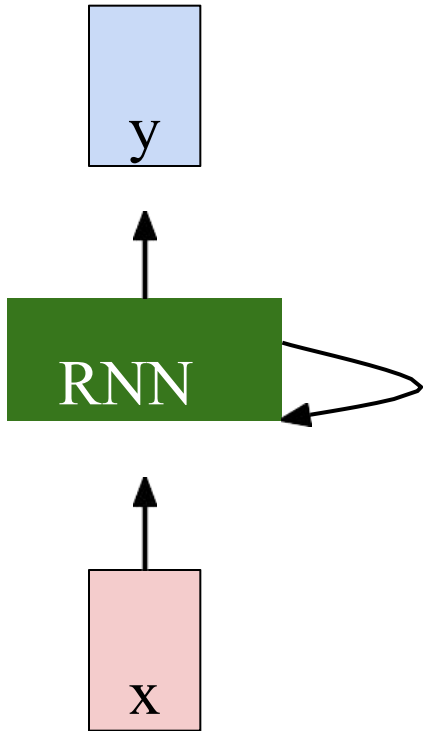
Notice: the same function and the same set of parameters are used at every time step.

# (Vanilla) Recurrent Neural Network

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

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

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

$$y_t = W_{hy} h_t$$

y

RNN

x

# RNN: Computational Graph

# RNN: Computational Graph

# RNN: Computational Graph

# RNN: Computational Graph

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

# RNN: Computational Graph: Many to Many

# RNN: Computational Graph: Many to Many

# RNN: Computational Graph: Many to Many

# Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

# Training a RNN I: BPTT

- Backpropagation through time (BPTT): The training algorithm for updating network weights to minimize error including time

# Remember BackPropagation

1. Present a training input pattern and propagate it through the network to get an output.

1. Compare the predicted outputs to the expected outputs and calculate the error.

2. Calculate the derivatives of the error with respect to the network weights.

3. Adjust the weights to minimize the error.

4. Repeat.

# BPTT I: Loss

$$L\left(\{x^{(1)},\ldots,x^{(\tau)}\},\{y^{(1)},\ldots,y^{(\tau)}\}\right)$$
$$=\sum_t L^{(t)}$$
$$=-\sum_t \log p_{\text{model}}\left(y^{(t)} \mid \{x^{(1)},\ldots,x^{(t)}\}\right),$$

The total loss for a given sequence x(t)

$$\frac{\partial L}{\partial L^{(t)}} = 1.$$

Our goal is to calculate the gradients of the error with respect to our parameters U, W and V and and then learn good parameters using Stochastic Gradient Descent.

**Ref:** **https://www.geeksforgeeks.org/ml-back-propagation-through-time/**

# Vanishing/ Exploding gradients

18

# Vanishing/ **Exploding** gradient ...

- During training gradients explode/vanish easily because of depth-in-time →
  **Exploding**/**<span style="color:red">Vanishing</span>** gradients!

- **Vanishing Gradient** occurs when the derivative or slope will get **smaller and smaller** as we go backward with every layer during backpropagation.

- **Exploding gradient** occurs when the derivatives or slope will get **larger and larger** as we go backward with every layer during backpropagation. <span style="color:red">This situation is the exact opposite of the vanishing gradients.</span>

# Vanishing/ **Exploding** gradient...

- A **vanishing** Gradient problem occurs with the sigmoid and tanh activation function because the derivatives of these units are between 0 to 0.25 and 0 to 1 respectively.

- Therefore, the updated weight values are small, and the new weight values are very similar to the old weight values.

- A **Exploding** Gradient problem happens because of weights, not because of the activation function.

- Due to high weight values, the derivatives will also higher so that the new weight varies a lot to the older weight, and the gradient will never converge.

# Question

- Do FNNs with several hidden layers suffer from vanishing gradient problem?

# Standard Solutions

- Proper initialization of Weight Matrix

- Regularization of outputs or Dropout

- Use of ReLU Activations as it's derivative is either 0 or 1

# Gated Units

# Standard RNN

24

# Long-Short Term Memory (LSTM)

- LSTM:
  - an improvement over RNN and used for modeling complex sequential data.

  - has three inputs ($C_{t-1}$, $h_{t-1}$, $X_t$) and three outputs($C_t$, $h_t$, $O_t$), where $O_t$ is the softmax of $h_t$.

  - **gating** mechanism: that controls the memoizing process.

  - **gates** are FC leyer followed an activation function:
    - forget gate
    - input gate
    - cell-state gate
    - output gate

- **cell state**: encode a kind of aggregation of data from all previous time-steps that have been processed,

- **hidden state:** is meant to encode a kind of characterization of the most recent time-step.



Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy

Schmidhuber, Jürgen, and Sepp Hochreiter. "Long short-term memory." Neural Comput 9, no. 8 (1997): 1735-1780.

# Long-Short Term Memory (LSTM)...

- LSTM:

  - consists of:

    - five activation functions (**sigmoid** or **tanh**), and

    - four math operations (3 multi, 1 add)



$$\phi(z) = \frac{1}{1+e^{-z}}$$

$$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Schmidhuber, Jürgen, and Sepp Hochreiter. "Long short-term memory." Neural Comput 9, no. 8 (1997): 1735-1780.

# Long-Short Term Memory (LSTM)...

**Forget Gate**



**Forget Gate:**

What part of memory to "forget" – zero means forget this bit

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

Concatenate

- Previous hidden-state and current input are concatenated pass through FC and a sigmoid activation function.

# Long-Short Term Memory (LSTM)...

**Input Gate**



What bits to insert into the next states

**Input Gate Layer**

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

**New contribution to cell state**

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

Classic neuron

What content to store into the next state

- Previous hidden state and current input are **concatenated**, and are passed through two parallel branches:
  - the first branch has a FC network, followed by a sigmoid,
  - second branch has a FC network, followed by a tanh function.

# Long-Short Term Memory (LSTM)...

**Update Cell State**



Next memory cell content –
mixture of not-forgotten part
of previous cell and insertion

**Update Cell State (memory):**

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

- The output vectors of two parallel branches are multiplied elementwise.
- The output after multiplication is added to the cell state after the cell state is updated with the forget gate.

# Long-Short Term Memory (LSTM)...

**Output Gate**



What part of cell to output

**Output Gate Layer**

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

**Output to next layer**

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

- The current input is concatenated with the hidden state and passed through the Fully connected network

- once appling the sigmoid, and this output is multiplied element-wise with cell state after it passes through a tanh function.

- The result of this multiplication is the current output/hidden state.

**Hint**:

- we have four FC networks at different places in the LSTM.

- each FC network has an input layer and an output layer.

- the **number of neurons** in each FC network input and output is dependent on the dimension of the input vector($m$) and cell state vector($n$).

- Since the input to all FC networks is a concatenated vector of current input and hidden state vector($h_{t-1} + X_t$). The input layer of all the FC networks has ($m+n$) neurons.

- Since the output of the FC network will be multiplied or added to the cell state element-wise, the output layer of FC network has a number of neutrons equal to the dimension of the cell state($n$).

## How to compute parameter LSTM?

**Hint**:

- we have four FC networks at different places in the LSTM.

- each FC network has an input layer and an output layer.

- the **number of neurons** in each FC network input and output is dependent on the dimension of the input vector($m$) and cell state vector($n$).

- Since the output of the FC network will be multiplied or added to the cell state element-wise, the output layer of FC network has a number of neutrons equal to the dimension of the cell state($n$).

- To calculate number of parameters. Each FC network has a parameter weight matrix of $[(m+n)*n]$ and a bias values of '$n$'.

- Total parameters at each FC network **($m*n$ + sqr($n$) + $n$)** and for four FC networks it will be **4*($m*n$ + sqr($n$) + $n$)**.

## Why do we need a cell state when we have a hidden state of LSTM?
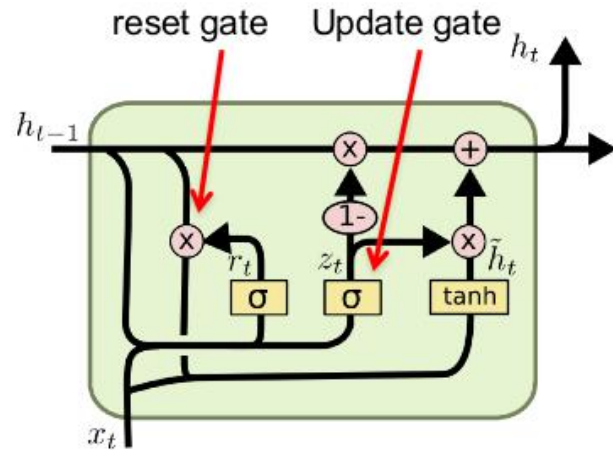
**Hint**: Vanishing/exploding gradient issues.

- if no cell-state: the hidden state feature vector will pass through the Neural network and gets to interact with weight matrices on the way.

- then this cell-state acts as a **highway without touching** any fully connected layers on the way, it is literally similar to the residual block in a resent architecture.

# GRU



LSTM

## GRU – gated recurrent unit
### (more compression)

reset gate    Update gate

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

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

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

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

It combines the forget and input into a single update gate. It also merges the cell state and hidden state. This is simpler than LSTM. There are many other variants too.

X,*: element-wise multiply

Reading assignment on GRU...?

# Open research Topics [CNN, RNN, CNN+RNN]

- Speech to text and via for low-resource languages [langues have no enough written documents]
- Visual question answering

- Image captioning
- Language identification from document images [e.g historical documents]
- Ethiopian sign language translation and recogntion
- Agricultural production quality grading and identification (e.g banana coffee, barley etc)
- Multi-modal data analysis( e.g taking the image and medical notes, prescription for disease diagnosis)
- Covid-19...impact prediction in a certain domain ( education, economy, health sectors...)
- Text detection and recognition in the wild

LSTM
CTC-decoder
P1 P2 P3 ... Pt
Input Image
Convolutional layers
BLSTM
Soft-max outputs
output text
Transcription layer

$(h_0,c_0)$ $(h_1,c_1)$ $(h_2,c_2)$ $(h_3,c_3)$ $(h_4,c_4)$ $(h_5,c_5)$ $(h_6,c_6)$ $(h_7,c_7)$
DECODER
<SOS>
<EOS>
$(h_0,c_0)$ $(h_1,c_1)$ $(h_2,c_2)$ $(h_3,c_3)$ $(h_4,c_4)$ $(h_5,c_5)$ $(h_6,c_6)$ $(h_7,c_7)$ $(h_8,c_8)$
ENCODER
learning is the process of gaining new understanding

$g_0$ $g_1$ ...
$p(g_0)$ $p(g_1)$ ... $p(g_7)$
CNN RNN RNN ... RNN
$y_0$ ... $y_{t-1}$

OUTPUT
ENCODER    internal states    DECODER
INPUT

fine <EOL>
am
W
How are you <EOL>
LSTM Encoder

LSTM Decoder
LSTM
LSTM
P1 P2 ... Pt
CTC Error calculation and Decoding
Input image
BLSTM
Soft-max output
Transcription

$h_{t-1}$ $h_t$ $h_{t+1}$
A A
$\sigma$ $\sigma$ tanh $\sigma$
tanh
$x_{t-1}$ $x_t$ $x_{t+1}$